

# Maximum softly-penalized likelihood for Bernoulli-response GLMMs

Philipp Sterzinger, Ioannis Kosmidis

Department of Statistics, University of Warwick, Coventry, UK  
philipp.sterzinger@warwick.ac.uk

WARWICK  
THE UNIVERSITY OF WARWICK

We provide a computationally stable estimation method for Bernoulli-response GLMMs that returns interior point estimates and has optimal asymptotic properties

## Bernoulli-response GLMMs

### Model

$$Y_{ij} | \mathbf{u}_i \sim \text{Bernoulli}(\mu_{ij}) \quad \text{with} \quad g(\mu_{ij}) = \eta_{ij} = \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \mathbf{z}_{ij}^\top \mathbf{u}_i$$

$$\mathbf{u}_i \sim \text{N}(\mathbf{0}_q, \boldsymbol{\Sigma}) \quad (i = 1, \dots, k; j = 1, \dots, n_i)$$

where  $\mu_{ij} = P(Y_{ij} = 1 | \mathbf{u}_i, \mathbf{x}_{ij}, \mathbf{z}_{ij})$ ,  $g: (0, 1) \rightarrow \mathbb{R}$  is a known link function, e.g.  $g(x) = \log(x/1-x)$ , and  $\mathbf{x}_{ij}, \mathbf{z}_{ij}$  are vectors of fixed and random effect covariates.

### Marginal likelihood about $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$

$$L(\boldsymbol{\beta}, \boldsymbol{\Sigma}) \propto \det(\boldsymbol{\Sigma})^{-k/2} \prod_{i=1}^k \int_{\mathbb{R}^q} \prod_{j=1}^{n_i} \mu_{ij}^{y_{ij}} (1 - \mu_{ij})^{1-y_{ij}} \exp \left\{ -\frac{\mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i}{2} \right\} d\mathbf{u}_i \quad (1)$$

### Maximum approximate likelihood (MAL) estimator

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta})$$

where  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\Sigma})$  and  $\ell(\boldsymbol{\theta})$  is the logarithm of an approximation to (1).

## Maximum softly penalized likelihood

### Maximum softly penalized approximate likelihood (MSPAL) estimator

$$\tilde{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \Theta} \{\ell(\boldsymbol{\theta}) + P(\boldsymbol{\theta})\}$$

where  $P(\boldsymbol{\theta})$  is a penalty that satisfies  $\lim_{r \rightarrow \infty} P(\boldsymbol{\theta}(r)) = -\infty$ , for any sequence of parameters  $\boldsymbol{\theta}(r)$  such that  $\lim_{r \rightarrow \infty} \boldsymbol{\theta}(r)$  lies on the boundary.

Then, if there is a point  $\boldsymbol{\theta}$  such that  $\ell(\boldsymbol{\theta}) + P(\boldsymbol{\theta}) > -\infty$ ,  $\tilde{\boldsymbol{\theta}}$  is in the interior of the parameter space.

### “Soft” penalty

To preserve the MAL asymptotics, we adapt the work of [3] and control  $\|\nabla P(\boldsymbol{\theta})\|$  in terms of the rate of information accumulation  $r_n$ , for which the observed information  $J(\boldsymbol{\theta}) = -\nabla \nabla^\top \ell(\boldsymbol{\theta})$  converges to a nonrandom  $\mathcal{O}(1)$  matrix  $I(\boldsymbol{\theta})$ , i.e.  $r_n^{-1} J(\boldsymbol{\theta}) \xrightarrow{p} I(\boldsymbol{\theta})$ .

### Consistency

A0 Both  $\ell(\boldsymbol{\theta}), \tilde{\ell}(\boldsymbol{\theta})$  are differentiable, with derivatives  $S(\boldsymbol{\theta}), \tilde{S}(\boldsymbol{\theta})$

A1  $\sup_{\boldsymbol{\theta} \in \Theta} \|r_n^{-1} S(\boldsymbol{\theta}) - S_0(\boldsymbol{\theta})\| \xrightarrow{p} 0$  for some deterministic function  $S_0(\boldsymbol{\theta})$

A2 For all  $\varepsilon > 0$ ,  $\inf_{\boldsymbol{\theta} \in \Theta: \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \geq \varepsilon} \|S_0(\boldsymbol{\theta})\| > 0 = \|S_0(\boldsymbol{\theta}_0)\|$

A3  $\hat{\boldsymbol{\theta}}$  and  $\tilde{\boldsymbol{\theta}}$  are roots of  $S(\boldsymbol{\theta}), \tilde{S}(\boldsymbol{\theta})$ , i.e.  $S(\hat{\boldsymbol{\theta}}) = \mathbf{0}$  and  $\tilde{S}(\tilde{\boldsymbol{\theta}}) = \mathbf{0}$

#### Theorem 1

Let  $\sup_{\boldsymbol{\theta} \in \Theta} \|\nabla P(\boldsymbol{\theta})\| = o_p(r_n)$ , and assume that A0-A3 hold. Then  $\tilde{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$ .

### Asymptotic Normality

A4 Both  $\ell(\boldsymbol{\theta}), \tilde{\ell}(\boldsymbol{\theta})$  are three times differentiable

A5  $\sup_{\boldsymbol{\theta} \in \Theta} \|\|r_n^{-1} J(\boldsymbol{\theta}) - I(\boldsymbol{\theta})\|\| \xrightarrow{p} 0$  and  $I(\boldsymbol{\theta})$  is continuous in  $\boldsymbol{\theta}$  in a neighbourhood around  $\boldsymbol{\theta}_0$

A6  $r_n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} \text{N}(0, I(\boldsymbol{\theta}_0)^{-1})$

A7  $\tilde{\boldsymbol{\theta}}$  is consistent for  $\boldsymbol{\theta}_0$

#### Theorem 2

Let  $\sup_{\boldsymbol{\theta} \in \Theta} \|\nabla P(\boldsymbol{\theta})\| = o_p(r_n^{1/2})$  and assume that conditions A3-A7 hold. Then

$$r_n^{1/2}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} \text{N}(0, I(\boldsymbol{\theta}_0)^{-1})$$

## Soft penalties for Bernoulli-response GLMMs

We propose to penalize the approximate log-likelihood  $\ell(\boldsymbol{\theta})$  by a composite penalty that penalizes the fixed effects  $\boldsymbol{\beta}$  and the Cholesky factor of the variance components, say  $\mathbf{L}$ , separately, i.e.

$$P(\boldsymbol{\theta}) = P^{FE}(\boldsymbol{\beta}) + P^{VC}(\mathbf{L})$$

### Fixed effects penalty

$$P^{FE}(\boldsymbol{\beta}) = \sqrt{p/n} \log \det(\mathbf{X}^\top \mathbf{W} \mathbf{X})$$

where  $\mathbf{X}$  is the  $n \times p$  matrix of all fixed effect covariates,  $\mathbf{W}$  is a diagonal matrix with entries  $W_{ii} = \mu_i(\boldsymbol{\beta})(1 - \mu_i(\boldsymbol{\beta}))$ ,  $\log(\mu_i(\boldsymbol{\beta})/(1 - \mu_i(\boldsymbol{\beta}))) = \mathbf{x}_i^\top \boldsymbol{\beta}$ .

[4] show that  $\lim_{r \rightarrow \infty} P^{FE}(\boldsymbol{\beta}(r)) = -\infty$  for any sequence of fixed effects such that  $\lim_{r \rightarrow \infty} \boldsymbol{\beta}(r)$  has infinite components, which guarantees finite MSPAL fixed effects estimates.

We show that  $\sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\nabla P^{FE}(\boldsymbol{\beta})\| = \mathcal{O}(1)$ , whenever  $\max_{i,j} |x_{ij}| = \mathcal{O}(n^{1/2})$  in line with Theorems 1 & 2.

### Variance components penalty

$$P^{VC}(\mathbf{L}) = \sum_{i=1}^q P^V(\log(l_{ii})) + \sum_{i < j} P^{VC}(l_{ij}), \quad P^V(x) = \sqrt{p/n} \begin{cases} -\frac{1}{2} \{x\}^2 & \text{if } |x| \leq 1, \\ -|x| + \frac{1}{2} & \text{otherwise} \end{cases}$$

where  $\mathbf{L}$  is the lower-triangular Cholesky factor of the variance components, so that  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^\top$ . This penalty gives estimated variance components matrix  $\tilde{\boldsymbol{\Sigma}} = \tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top$  that are finite, nonsingular and without perfect estimated correlation, i.e. for all  $i \neq j$ ,  $\left| \frac{\tilde{\Sigma}_{ij}}{\sqrt{\tilde{\Sigma}_{ii}\tilde{\Sigma}_{jj}}} \right| < 1$ .

## Motivating Example

### Data

Table 1: Culcita data from the worked examples of [1]

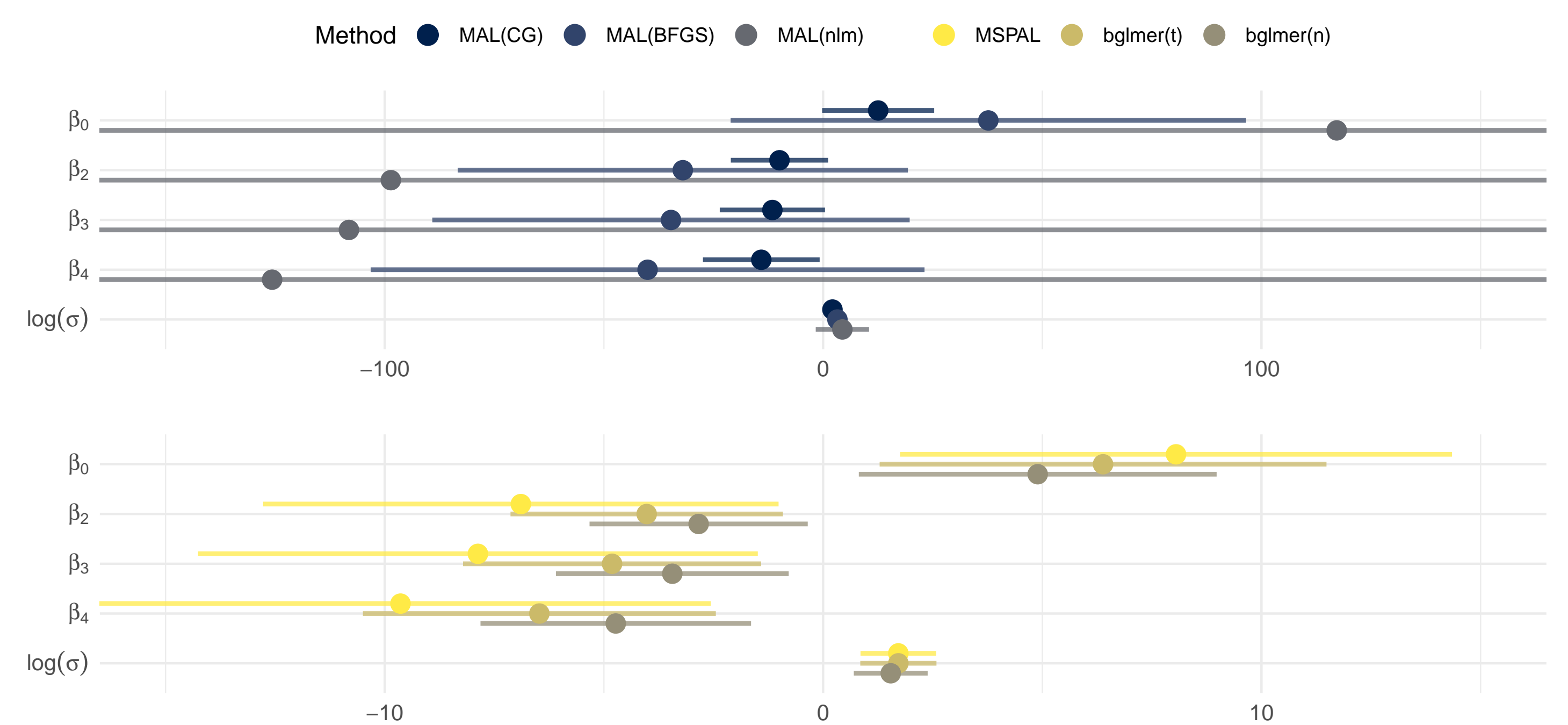
Treatment	Block									
	1	2	3	4	5	6	7	8	9	10
none	0,1	1,1	1,1	1,1	1,1	1,1	1,1	1,1	1,1	1,0
crabs	0,0	0,0	0,0	0,0	1,1	1,1	1,1	1,1	1,1	1,1
shrimp	0,0	0,0	0,0	0,0	0,1	1,1	1,1	1,1	1,1	1,1
both	0,0	0,0	0,0	0,0	0,0	0,1	1,1	1,1	1,1	1,1

The complete randomized block design (four treatments, ten temporal blocks, two replications) data records coral-eating sea stars (Culcita) attacking coral harbouring different protective symbionts (crabs, shrimp).

### Model

Upon removal of the atypical observation in block 10 without predation and symbionts, we associate predation to treatment effects using a Bernoulli-response GLMM with logistic link and a random intercept per block, so that  $\eta_{ij} = \beta_0 + u_i + \beta_j$ ,  $u_i \sim \text{N}(0, \sigma^2)$  for  $i = 1, \dots, 10, j = 1, \dots, 4$ .  $\beta_1$  is set to zero, making treatment “none” the reference category.

### Infinite & disparate estimates



We estimate  $\boldsymbol{\beta}, \log \sigma$  by MAL, with a 200-point Gauss-Hermite quadrature approximation to the log-likelihood and three different optimization routines from the R `optimx` package, and our proposed MSPAL estimator and `bgfmer` [2], with a 100-point adaptive Gauss-Hermite quadrature approximation to the log-likelihood.

The MAL estimates and asymptotic 95% CIs about the model parameters are highly dissimilar and effectively infinite but, due to different stopping criteria of the optimizers, may not be perceived as such. MSPAL and `bgfmer` estimates are finite, with `bgfmer` exhibiting stronger shrinkage.

## Simulation study



Simulation results for the MAL, MSPAL and the `bgfmer` routine of [2] from 10000 independent samples of responses at the MAL estimates when all data points are used.

## References

- [1] B. M. Bolker, “Linear and generalized linear mixed models,” in *Ecological Statistics*, Oxford University Press, 2015, pp. 309–333. DOI: 10.1093/acprof:oso/9780199672547.003.0014.
- [2] Y. Chung, S. Rabe-Hesketh, et al., “A nondegenerate penalized likelihood estimator for variance parameters in multilevel models,” *Psychometrika*, vol. 78, no. 4, pp. 685–709, 2013.
- [3] H. Ogden, “On asymptotic validity of naive inference with an approximate likelihood,” *Biometrika*, vol. 104, no. 1, pp. 153–164, 2017.
- [4] I. Kosmidis and D. Firth, “Jeffreys-prior penalty, finiteness and shrinkage in binomial-response generalized linear models,” *Biometrika*, vol. 108, no. 1, pp. 71–82, 2021.