

Accurate Frequentist Generalised Linear Mixed Model Analysis via Expectation Propagation

by **James Yu**

Thesis submitted in fulfilment of the requirements for the degree of

Doctor of Philosophy in Mathematics

under the supervision of **Prof. Matt Wand** and **Dr. Shev MacNamara**

University of Technology Sydney

Faculty of Science

April 22, 2021

Certificate of original authorship

I, **James Yu** declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy in Mathematics, in the School of Mathematical and Physical Sciences at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise reference or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of the requirements for a degree at any other academic institution except as fully acknowledged within the text. This thesis is the result of a Collaborative Doctoral Research Degree program with the Australian Red Cross Blood Service.

This research is supported by an Australian Government Research Training Program.

Production Note:

Signature: Signature removed prior to publication.

Date: April 22, 2021

Acknowledgements

This thesis would not have been possible without the support of a number of individuals.

Many thanks to my mentor and supervisor Matt Wand. Thank you for your patience and kindness. This achievement would not have been possible without you. I could not have asked for a better supervisor.

I would also like to thank Stephen Wright. Thank you for your care and for the time you spent with me. I really appreciate it.

To my colleagues and friends at UTS, thank you for your help, constant support, for making me laugh and for listening to me.

Finally to Mum, Dad and Mario, I am eternally grateful to you for giving me the opportunities and experiences that have made me who I am. I dedicate this milestone to you.

List of papers/publications

The following list of paper and publications awards relate to the work presented in this thesis:

- Hall, P., Johnstone, I. M., Ormerod, J. T., Wand, M. P., & Yu, J. C. F. (2020). Fast and Accurate Binary Response Mixed Model Analysis via Expectation Propagation, *Journal of the American Statistical Association*, 115:532, 1902–1916, DOI: 10.1080/01621459.2019.1665529.
- Wand, M. P. & Yu, J. C. F. (2020). *glmmEP: Generalized Linear Mixed Model Analysis via Expectation Propagation (Version 1.0-3.1)*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/glmmEP/index.html>
- Yu, J. C. F. (2019), “Fast and Accurate Frequentist Generalised Linear Mixed Model Analysis via Expectation Propagation”, (2019) Enabling Algorithms Theme Symposium, University of Technology Sydney, Sydney, 13-14 June.
- Wand, M. P., & Yu, J. C. F. (2021) Density Estimation via Bayesian Inference Engines.
- First Place, 20th Annual J.B. Douglass Award, New South Wales Branch of Statistical Society of Australia, (2019).

Notation

In this chapter, important notations are introduced that we refer to throughout this thesis.

Acronyms

Table 1: Table with acronyms used in the thesis with their meanings.

Acronym	Meaning
AGHQ	Adaptive Gauss-Hermite quadrature
BFGS	Broyden-Fletcher-Goldfarb-Shanno
BP	Best predictor
CDF	Cumulative density function
DAG	Directed acyclic graph
DC	Data cloning
EP	Expectation propagation
GHQ	Gauss-Hermite quadrature
GLMM	Generalised linear mixed models
KL	Kullback-Leiber
MCMC	Markov chain Monte Carlo
NM	Nelder-Mead
PDF	Probability density function
PQL	Penalised quasi-likelihood

Contents

1 Introduction and background	2
1.1 Introduction	2
1.2 Thesis aim	4
1.3 Matrix storage and notation	5
1.4 Notation for spaces	7
1.5 Exponential family theory and distributions	7
1.5.1 Exponential family theory	7
1.5.2 Probability distribution	9
1.5.2.1 Univariate normal distribution	9
1.5.2.2 Multivariate normal distribution	10
1.5.2.3 Bernoulli distribution	11
1.5.2.4 Poisson distribution	12
1.5.2.5 Negative binomial distribution	13
1.6 Graph theory	13
1.6.1 Directed acyclic graphs	14
1.6.2 Factor graphs	15
1.7 Multilevel datasets	15
1.8 Generalised linear mixed models	19
1.8.1 Binary response models	22
1.8.2 Count response models	22
1.9 Maximum likelihood	23
1.9.1 Likelihood functions	23
1.9.2 Maximum likelihood	25
1.10 Best prediction	25
1.11 Current approximation methods	26
1.11.1 Laplace approximation	27
1.11.2 Gauss Hermite quadrature	28
1.11.3 Other methods	29
1.12 Expectation propagation	30
1.12.1 Kullback Leiber divergence and projection	31
1.12.2 Mean field approximations	33
1.12.3 Expectation propagation	34
1.12.4 Message passing	37
1.13 Thesis structure	39

2	Expectation propagation for the simplest one level probit mixed model	40
2.1	Traditional quadrature likelihood approximation	41
2.2	Expectation propagation likelihood approximation	42
2.2.1	Projection onto the unnormalised normal family	45
2.2.2	Message passing formulation	47
2.2.3	Starting values for Algorithm 3	54
2.3	Evaluation of the estimates	55
2.4	Computing point estimates and confidence intervals	56
2.4.1	Confidence interval estimation	57
2.4.2	Derivative approximation	58
2.4.2.1	Expectation propagation analytical approach I	61
2.4.2.2	Expectation propagation analytical approach II	63
2.5	Best predictor	70
2.6	Simulation Study	71
2.7	Appendix	75
2.7.1	Proof of Result 3	75
2.7.2	Proof of Result 5	76
2.7.2.1	Proof of Lemma 1	77
2.7.3	Proof of Result 7	80
2.7.4	Details on finding the inverse function gradient map log-partition function for each s	81
2.7.4.1	The $s = 0$ case	81
2.7.4.2	The $s = 2$ case	82
2.7.4.3	The $s = 4$ case	82
3	Expectation propagation for general one level probit mixed models	83
3.1	Expectation propagation likelihood approximation	84
3.1.1	Projection onto the unnormalised multivariate normal family	88
3.1.2	Message passing formulation	89
3.1.3	Starting values for Algorithm 6	95
3.2	Computation of point estimates and confidence intervals	96
3.3	Best predictor	99
3.4	Simulation study	100
3.4.1	Comparison of maximum likelihood estimates for univariate ran- dom effects	100
3.4.2	Maximum likelihood estimates for bivariate random effects	104
3.4.3	Maximum likelihood estimates for trivariate random effects	106
3.5	Appendix	108
3.5.1	Proof of Definition 13	108
3.5.1.1	Proof of Lemma 2	108
3.5.1.2	Proof of Lemma 3	109
3.5.2	Proof of Result 12	111
4	Expectation propagation for one level logistic mixed models	113
4.1	The simplest logistic mixed model	114
4.1.1	Traditional quadrature likelihood approximation	115
4.1.2	Expectation propagation likelihood approximation	115

4.1.2.1	Projection onto the unnormalised normal family	118
4.1.2.2	Message passing formulation	119
4.1.2.3	Starting values for the univariate logistic case	120
4.1.3	Evaluation of the estimates	121
4.1.4	Best predictor	123
4.2	General logistic mixed models	123
4.2.1	Expectation propagation likelihood approximation	124
4.2.1.1	Projection onto the unnormalised multivariate normal family	127
4.2.1.2	Message passing formulation	130
4.2.1.3	Starting values for the multivariate logistic case	131
4.2.2	Simulation study	132
4.2.2.1	Comparison of maximum likelihood estimates for univariate random effects	132
4.2.2.2	Maximum likelihood estimates for bivariate random effects	134
4.3	Appendix	137
4.3.1	Proof of Result 16	137
4.3.2	Proof of Definition 17	138
4.3.3	Proof of Definition 18	140
5	Expectation propagation for one level count response mixed models	142
5.1	The simplest Poisson mixed model	143
5.1.1	Traditional quadrature likelihood approximation	144
5.1.2	Expectation propagation likelihood approximation	144
5.1.2.1	Projection onto the unnormalised normal family	146
5.1.2.2	Message passing formulation	147
5.1.2.3	Starting values for the Poisson case	148
5.1.3	Evaluation of the estimates	150
5.1.4	Best predictor	151
5.2	General Poisson mixed models	151
5.2.1	Expectation propagation likelihood approximation	152
5.2.1.1	Projection onto the unnormalised multivariate normal family	154
5.2.1.2	Message passing formulation	155
5.2.1.3	Starting values for the Poisson case	156
5.2.2	Simulation study	157
5.3	The simplest negative binomial models	160
5.3.1	Traditional quadrature likelihood approximation	161
5.3.2	Expectation propagation likelihood approximation	161
5.3.2.1	Projection onto the unnormalised normal family	163
5.3.2.2	Message passing formulation	165
5.3.2.3	Starting values for the negative binomial case	166
5.3.3	Evaluation of the estimates	167
5.3.4	Computation of point estimates and confidence intervals	168
5.3.4.1	Derivative approximation	169
5.4	General negative binomial model	169
5.4.1	Expectation propagation likelihood approximation	170

5.4.1.1	Projection onto the unnormalised multivariate normal family	172
5.4.1.2	Message passing formulation	174
5.4.1.3	Starting values for the negative binomial case	175
5.4.2	Computation of point estimates and confidence intervals	176
5.4.3	Simulation study	177
5.5	Varying dispersion negative binomial model	179
5.6	Appendix	181
5.6.1	Proof of Definition 21	181
5.6.2	Proof of Definition 23	182
5.6.3	Proof of Definition 25	184
5.6.4	Proof of Definition 27	185
6	Expectation propagation for two level and crossed random effects probit models	187
6.1	The general probit crossed mixed model	187
6.1.1	Expectation propagation likelihood approximation	188
6.1.1.1	Message passing formulation	189
6.1.2	Computation of point estimates and confidence intervals	196
6.1.3	Simulation study	199
6.1.3.1	Comparison with MCMC and Laplace approximation	
	maximum likelihood for crossed random effects	199
6.2	The general probit two level mixed model	200
6.2.1	Expectation propagation likelihood approximation	202
6.2.1.1	Message passing formulation	202
7	Applications of expectation propagation for one level probit mixed models	210
7.1	R package “glmmEP”	210
7.2	Modelling immunisation of Guatemalan children	211
7.3	Modelling donor attendance of the Australian Red Cross Blood Service	216
7.3.1	Data cleaning	220
7.3.2	Modelling continuous variables	222
7.3.3	Initial model	223
7.3.4	Second model	226
7.3.4.1	Results of model fit	227
7.4	Appendix	233
8	Discussion and conclusion	235
9	References	239

Abstract

Generalised linear mixed models are a particularly powerful and well established statistical tool. Unlike linear mixed models, where the integrals arising in likelihood functions can be expressed in closed form, the likelihood functions expressed in generalised linear mixed models do not follow tractable solutions. Methods such as Gauss-Hermite quadrature and Laplace approximation are the standard approaches to overcome these integrals. Although Gauss-Hermite quadrature is accurate it is also slow, rendering it unsuitable for analyses with more than two or three random effects. Laplace approximations are the most feasible solution, however the approximate inference they provide in binary models is well known to be inaccurate. A less common approach is to use Bayesian ideas such as data cloning, however they involve a number of technicalities and as such are difficult to implement. Although expectation propagation is generally used in Bayesian settings, in this thesis we introduce a novel approach where we use it as frequentists to achieve high accuracy results with minimal computational cost for inference on generalised linear mixed models. We show our methodology can be used to solve one level probit models without the need for quadrature, providing consistent and accurate results. We explain how using quadrature we can also extend our method to logistic, Poisson and negative-binomial models. Additionally we show how these models can be extended to two level models and crossed random effects models for the probit case. Finally we present applications of our methodology on two real datasets, both with different technical challenges.

Chapter 1

Introduction and background

This chapter provides an introduction to the thesis as well as background information and theory required for the upcoming chapters. None of the work presented in this chapter is novel.

1.1 Introduction

In reality analysts do not always receive data that conforms to the properties required for the implementation of classical statistical procedures. Data often contains non-normally distributed response variables and heterogeneous variance. These variance structures may be explained by multilevel or hierarchical structures, where the heterogeneity stems from observations nested within larger groups of experimental units (Steenbergen & Jones, 2002^[62]). As an example consider a dataset of results from an experiment that is repeated 100 times by 5 different people (there are a total of 500 replicates). Since each person will conduct the experiment slightly differently, modelling the data requires accounting for variability both within and across persons. This type of dataset is common in biological, medical and social sciences, and is becoming increasingly prevalent in other fields. Section 1.2 of Gelman & Hill (2007^[19]) and the R package “mlmRev” (Bates, Maechler & Bolker, 2019^[4]) provide real examples. Although it may be tempting to coerce data into classical statistical frameworks via tricks such as data transformations or ignoring the effect of data structures, doing so may violate key statistical assumptions and limit the scope of inference gained (Bolker, et al., 2009^[8]). Subsequently, a more suitable approach of analysis is required.

Generalised linear mixed models (GLMMs) are a particularly powerful and well

established statistical tool. Like generalised linear models GLMMs can model non-normal responses, however they additionally account for nested data structures. GLMMs do this by allowing a population intercept and slope (known as fixed effects) as well as a unique intercept and slope for each experimental unit (known as random effects). These are the attributes that allow GLMMs to handle combinations of data structures and response distributions.

Unlike linear mixed models where the integrals required to solve likelihood functions can be obtained analytically, they do not have tractable solutions in GLMMs. Additionally, as the number of random effects increases they become more computationally intensive. In the frequentist setting, methods such as Gauss-Hermite quadrature (GHQ), penalised quasi-likelihood (PQL) and Laplace approximation are the standard approaches to overcome these integrals. PQL is one of the simplest and most widely used methods, available in a variety of computing software. However, it deals poorly with binary data where the standard deviations of random effects are large, often producing biased parameter estimates (Bolker, et al., 2009^[8]). Additionally, rather than approximate a true likelihood, PQL approximates a quasi-likelihood, making likelihood based inference infeasible. Both GHQ and Laplace approximation are more accurate than PQL, with GHQ the most accurate of the three. However, the accuracy of GHQ comes at the price of speed. As the number of random effects increases GHQ slows considerably, rendering it unsuitable for analyses with more than two random effects. While Laplace approximations provide the most feasible solution, the approximate inference it provides in binary models is inaccurate (McCulloch, et al., 2008^[41]), particularly when the number of observations per grouping unit are low.

Lele, et al. (2007)^[35] present a reformulation of the typically Bayesian method Markov chain Monte Carlo, known as data cloning, which allows for the calculation of maximum likelihood estimates and confidence intervals. Although the proposed method is able to cover random effects, it also involves several difficult technical details regarding its fitting and implementation. Similar issues exist for the R package “MCMCglmm” (Hadfield, 2017^[25]).

Perhaps a lesser known method, expectation propagation (EP) (Minka, 2001^[43]) has underpinnings originating from computer science. A variety of software such as Stan (Stan Development Team, 2017^[24]) and Infer.Net (Minka, et al., 2014^[45]) exist and facilitate the implementation of such deterministic algorithms for fast inference. Although the speed of computation is increased over more traditional methods, algebraic overheads mean utilisation of this methodology is cumbersome. Minka (2005)^[44] partially solves this issue by developing a message passing approach to EP, which allows the

algebra to be broken into fragments. Although the initial algebraic overhead is still large, once calculated, fragments can be reused via graphical structures without any additional algebra. Further numerical studies have shown EP to be more accurate than variational alternatives such as mean field variational Bayes and easily faster than Markov chain Monte Carlo. Specifically, Kim & Wand (2016)^[31] provide computational studies which suggest EP becomes more accurate than mean field variational Bayes as sample size increases (Kim & Wand, 2017^[32]).

Kim & Wand (2016)^[31] demonstrate the form of EP for the simple statistical problem of Bayesian inference from independent and identically distributed observations from a normal distribution. This paper also provides the means to implement a message passing approach using factor graphs, which facilitates extension to larger models with minimal algebraic overheads. Kim & Wand (2017)^[32] provide the structure of algorithms required for implementing EP in GLMM settings.

1.2 Thesis aim

Although many methods for GLMM analysis in a frequentist setting exist, they provide either poor accuracy or bad computational performance. Additionally, current applications of EP in statistics are primarily limited to Bayesian settings. This thesis aims to develop novel methodology for GLMM analysis in a frequentist setting which utilises EP for consistent and accurate inference, particularly in the case when low numbers of observations per group variable occur. We aim to build on Kim & Wand (2017)^[32] by utilising message passing to streamline the computations from the factor graph. We aim to provide frameworks for one level probit, logistic, Poisson and negative binomial models, as well as two level and crossed random effects probit models. Finally we aim to create an R package and make the methodology available, with a demonstration on a dataset provided by the Australian Red Cross Blood Service.

1.3 Matrix storage and notation

For a vector \mathbf{a} of length d , we denote the elements from i to j by $\mathbf{a}_{i:j}$. For clarity, consider

$$\mathbf{a} = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}, \quad \text{then} \quad \mathbf{a}_{1:2} \equiv \begin{bmatrix} a_0 \\ a_1 \end{bmatrix}.$$

To combat issues arising from storage of large matrices we now introduce concepts and notation analogous to Magnus & Neudecker (1999).^[38] For a $d \times d$ matrix \mathbf{A} , $\text{vec}(\mathbf{A})$ returns a $d^2 \times 1$ vector, where the columns of \mathbf{A} are stacked underneath each other in order from left to right. $\text{vech}(\mathbf{A})$ returns a $\frac{1}{2}d(d+1) \times 1$ vector, where the entries above the diagonal are removed and the remaining entries are stacked by column from left to right. $\text{vecbd}(\mathbf{A})$ returns a $\frac{1}{2}d(d-1) \times 1$ vector, where the entries below the diagonal are stacked by column from left to right. For example, if

$$\mathbf{A} \equiv \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix},$$

then

$$\text{vec}(\mathbf{A}) \equiv \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \\ a_{12} \\ a_{22} \\ a_{32} \\ a_{13} \\ a_{23} \\ a_{33} \end{bmatrix}, \quad \text{vech}(\mathbf{A}) \equiv \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \\ a_{22} \\ a_{32} \\ a_{33} \end{bmatrix} \quad \text{and} \quad \text{vecbd}(\mathbf{A}) \equiv \begin{bmatrix} a_{21} \\ a_{31} \\ a_{32} \end{bmatrix}. \quad (1.1)$$

For a vector \mathbf{a} of length d^2 , $\text{vec}^{-1}(\mathbf{a})$ returns a $d \times d$ matrix. The i th column corresponds to the entries

$$\mathbf{a}_{(d(i-1)+1):(di)}.$$

For a vector \mathbf{a}^* of length $\frac{1}{2}d(d+1)$, $\text{vech}^{-1}(\mathbf{a}^*)$ returns a symmetric $d \times d$ matrix. The i th lower diagonal column and i th upper diagonal row correspond to the entries

$$\mathbf{a}^*_{(d(i-1)-\sum_{j=1}^{i-1}(j-1)+1):(di-\sum_{j=1}^i(j-1))}.$$

As an example, where $\mathbf{a} = \text{vec}(\mathbf{A})$ and $\mathbf{a}^* = \text{vech}(\mathbf{A})$,

$$\text{vec}^{-1}(\mathbf{a}) \equiv \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}, \quad \text{vech}^{-1}(\mathbf{a}^*) \equiv \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{21} & a_{22} & a_{32} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

Note that when \mathbf{A} is symmetric

$$\text{vech}^{-1}(\text{vech}(\mathbf{A})) = \mathbf{A}.$$

We denote the transpose of \mathbf{A} by \mathbf{A}^\top , and assuming \mathbf{A} is invertible we denote its inverse by \mathbf{A}^{-1} . The 3×3 diagonal matrix formed setting the upper and lower off diagonal entries of \mathbf{A} to 0 is given by

$$\text{diag}(\mathbf{A}) \equiv \begin{bmatrix} a_{11} & 0 & 0 \\ 0 & a_{22} & 0 \\ 0 & 0 & a_{33} \end{bmatrix},$$

where \mathbf{A} is defined in equation (1.1). The $d \times d$ identity matrix is denoted as \mathbf{I}_d , the $d \times 1$ vector with all entries equal to zero as $\mathbf{0}_d$, and the $d_1 \times d_2$ matrix with all entries equal to zero as $\mathbf{0}_{d_1 \times d_2}$.

The duplication matrix of order d is the unique $d^2 \times \frac{1}{2}d(d+1)$ matrix \mathbf{D}_d of zeros and ones such that

$$\mathbf{D}_d \text{vech}(\mathbf{A}) = \text{vec}(\mathbf{A}),$$

where \mathbf{A} is symmetric. The R function “`duplication.matrix()`” in the “`matrixcalc`” package (Novomestky, 2020^[48]) allows for easy calculation of the duplication matrix for any matrix \mathbf{A} . The Moore-Penrose inverse of \mathbf{D}_d is

$$\mathbf{D}_d^+ \equiv (\mathbf{D}_d^\top \mathbf{D}_d)^{-1} \mathbf{D}_d^\top.$$

The norm of a vector \mathbf{a} is defined as

$$\|\mathbf{a}\| = \sqrt{\mathbf{a}^\top \mathbf{a}}.$$

and for the same vector we define

$$\mathbf{a}^{\otimes k} \equiv \begin{cases} 1 & \text{for } k = 0 \\ \mathbf{a} & \text{for } k = 1 \\ \mathbf{a}\mathbf{a}^\top & \text{for } k = 2 \end{cases}. \quad (1.2)$$

1.4 Notation for spaces

We denote spaces as follows:

\mathbb{R} Set of real numbers.

\mathbb{R}^d Real coordinate space of d dimensions.

$\mathbb{R}_{\geq 0}$ Coordinate line with all real positive numbers (real numbers greater than or equal to zero).

$\mathbb{R}_{> 0}$ Coordinate line with all strictly positive real numbers (real numbers greater than zero).

$\mathbb{Z}_{\geq 0}$ Coordinate line with all positive integers (integers greater than or equal to zero).

1.5 Exponential family theory and distributions

1.5.1 Exponential family theory

The exponential family is a set of probability distributions which can be written in a specified parameteric form, where different distributions arise from varying parameter values.

Definition 1. *Given a natural parameter vector $\boldsymbol{\eta}$, a vector \boldsymbol{x} is from an exponential family of probability distributions if its probability distribution conditional on $\boldsymbol{\eta}$ can be expressed in the following form:*

$$f_{\boldsymbol{X}}(\boldsymbol{x}; \boldsymbol{\eta}) = \exp(\boldsymbol{\eta}^\top \boldsymbol{T}(\boldsymbol{x}) - A(\boldsymbol{\eta}))h(\boldsymbol{x}), \quad (1.3)$$

where $\boldsymbol{\eta}$ is the natural parameter vector, $\boldsymbol{T}(\boldsymbol{x})$ is the sufficient statistic, $A(\boldsymbol{\eta})$ is the log-partition function and $h(\boldsymbol{x})$ is the base measure.

Note that Definition [1](#) can be simplified for single parameter and univariate distributions.

We now clarify the previously arbitrary definitions of $\boldsymbol{\eta}$, $\boldsymbol{T}(\boldsymbol{x})$ and $A(\boldsymbol{\eta})$:

- $\boldsymbol{T}(\boldsymbol{x})$ is the sufficient statistic of the distribution, which is a function that given data \boldsymbol{x} provides all the information required to describe the posterior distribution of the natural parameters $\boldsymbol{\eta}$.
- $\boldsymbol{\eta}$ are the parameters of each distribution written in their natural form. Canonical link functions take distribution parameters as an argument and return natural parameters. The natural parameter space H is the convex space given by the set of values

$$H = \{\boldsymbol{\eta} : A(\boldsymbol{\eta}) < \infty\}.$$

- $A(\boldsymbol{\eta})$ is the log of the normalising factor

$$A(\boldsymbol{\eta}) = \log \left\{ \int h(\boldsymbol{x}) \exp(\boldsymbol{\eta}^\top \boldsymbol{T}(\boldsymbol{x})) d\boldsymbol{x} \right\},$$

which ensures that $f_{\boldsymbol{X}}(\boldsymbol{x}; \boldsymbol{\eta})$ is a probability density function. Additionally, the moments of the sufficient statistic $\boldsymbol{T}(\boldsymbol{x})$ can be derived by differentiating $A(\boldsymbol{\eta})$. More explicitly, the mean and variance of the sufficient statistic are respectively given by the first and second derivatives of $A(\boldsymbol{\eta})$,

$$E(\boldsymbol{T}(\boldsymbol{x})) = \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^\top = (\nabla A)(\boldsymbol{\eta}) \quad \text{and} \quad \text{Cov}(\boldsymbol{T}(\boldsymbol{x})) = \mathbf{D}_{\boldsymbol{\eta}} (\mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^\top) = \mathbf{H}(A(\boldsymbol{\eta})) \quad (1.4)$$

where $\text{Cov}(\boldsymbol{T}(\boldsymbol{x}))$ is the covariance matrix of $\boldsymbol{T}(\boldsymbol{x})$, $\mathbf{D}_{\boldsymbol{x}} f(\boldsymbol{x})$ is a $p \times d$ matrix with (i, j) entry equal to $df(\boldsymbol{x})_i/d\boldsymbol{x}_j$ for a \mathbb{R}^p -valued function f with argument $\in \mathbb{R}^d$, and $(\nabla A)(\boldsymbol{\eta})$ is a one-to-one function that returns a column vector of partial derivatives of $A(\boldsymbol{\eta})$ with respect to the entries of $\boldsymbol{\eta}$. A summary of $A(\boldsymbol{\eta})$ and $(\nabla A)(\boldsymbol{\eta})$ for exponential families is provided in Section 3.5 of Wainwright & Jordan (2008).^{[66](#)}

Exponential families have extremely useful properties in statistical analysis. As such, they have been used in classical statistics for decades and more recently in machine learning settings. A major benefit of exponential families is conjugate priors, where posterior distributions are in the same probability distribution family as the prior probability distribution. This helps reduce computational complexity for Bayesian inference.

1.5.2 Probability distribution

We now present the distributions used in this thesis and their exponential family parameterisation.

1.5.2.1 Univariate normal distribution

Definition 2. A scalar random variable x is from a univariate normal distribution with mean μ and variance $\sigma^2 > 0$ when its density function is

$$f(x) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right),$$

which is written as $x \sim N(\mu, \sigma^2)$.

Definition 3. A scalar random variable x is from a standard univariate normal distribution with $\mu = 0$ and $\sigma^2 = 1$ when its density function is

$$\phi(x) = (2\pi)^{-1/2} \exp\left(-\frac{x^2}{2}\right),$$

which is written as $x \sim N(0, 1)$, and cumulative distribution function (CDF)

$$\Phi(x) \equiv \int_{-\infty}^x \phi(t) dt.$$

We set

$$\zeta(x) = \log(2\Phi(x)),$$

with first and second derivatives respectively

$$\zeta'(x) = \Phi(x)/\phi(x) \quad \text{and} \quad \zeta''(x) = -\zeta'(x)(x + \zeta'(x)).$$

The sufficient statistic and base measure for the univariate normal distribution are respectively defined as

$$\mathbf{T}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad \text{and} \quad h(x) = (2\pi)^{-1/2}.$$

The natural parameter vector and inverse mapping are respectively

$$\boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} \mu/\sigma^2 \\ -1/(2\sigma^2) \end{bmatrix} \quad \text{and} \quad \begin{aligned} \mu &= -\eta_1/(2\eta_2) \\ \sigma^2 &= -1/(2\eta_2) \end{aligned}. \quad (1.5)$$

The log-partition function is

$$A(\boldsymbol{\eta}) = -\eta_1^2/(4\eta_2) - \frac{1}{2} \log(-2\eta_2).$$

The domain of both $A(\boldsymbol{\eta})$ and $(\nabla A)(\boldsymbol{\eta})$ is $H = (\eta_1, \eta_2) : \eta_1 \in \mathbb{R}, \eta_2 \in \mathbb{R}_{>0}$. The column vector of partial derivatives of $A(\boldsymbol{\eta})$ is given by

$$(\nabla A)(\boldsymbol{\eta}) = \begin{bmatrix} -\eta_1/(2\eta_2) \\ (\eta_1^2 - 2\eta_2)/(4\eta_2^2) \end{bmatrix}.$$

1.5.2.2 Multivariate normal distribution

Definition 4. A $d \times 1$ vector of random variables \mathbf{x} is from the multivariate normal distribution with $d \times 1$ mean vector $\boldsymbol{\mu}$ and $d \times d$ positive definite covariance matrix variance $\boldsymbol{\Sigma}$ when its density function is

$$f(\mathbf{x}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

which is written as $\mathbf{x} \sim \mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Definition 5. A $d \times 1$ vector of random variables \mathbf{x} is from the standard multivariate normal distribution with $d \times 1$ mean vector $\mathbf{0}_d$ and $d \times d$ positive definite covariance matrix variance \mathbf{I}_d when its density function is

$$\phi_{\mathbf{I}}(\mathbf{x}) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2}\mathbf{x}^\top \mathbf{x}\right),$$

which is written as $\mathbf{x} \sim \mathbf{N}(\mathbf{0}_d, \mathbf{I}_d)$.

The sufficient statistic and base measure for the multivariate normal distribution are respectively defined as

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix} \quad \text{and} \quad h(\mathbf{x}) = (2\pi)^{-d/2}.$$

The natural parameter vector and inverse mapping are respectively

$$\boldsymbol{\eta} \equiv \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\mathbf{D}_d^\top \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} \quad \text{and} \quad \begin{matrix} \boldsymbol{\mu} = -\frac{1}{2}\mathbf{H}_2^{-1}\boldsymbol{\eta}_1 \\ \boldsymbol{\Sigma} = -\frac{1}{2}\mathbf{H}_2^{-1} \end{matrix}, \quad (1.6)$$

where

$$\boldsymbol{\eta}_2 \equiv \mathbf{D}_d^\top \text{vec}(\mathbf{H}_2) \quad \text{and} \quad \mathbf{H}_2 = \text{vec}^{-1}((\mathbf{D}_d^+)^\top \boldsymbol{\eta}_2). \quad (1.7)$$

The log-partition function is

$$A(\boldsymbol{\eta}) = -\frac{1}{4}\boldsymbol{\eta}_1^\top \mathbf{H}_2^{-1}\boldsymbol{\eta}_1 - \frac{1}{2} \log | -2\mathbf{H}_2|,$$

The domain of both $A(\boldsymbol{\eta})$ and $\nabla A(\boldsymbol{\eta})$ is $H = (\boldsymbol{\eta}_1, \boldsymbol{\eta}_2) : \boldsymbol{\eta}_1 \in \mathbb{R}^d, \boldsymbol{\eta}_2 \in \mathbb{R}_{>0}^{\frac{1}{2}d(d+1)}$. The column vector of partial derivatives of $A(\boldsymbol{\eta})$ is

$$(\nabla A)(\boldsymbol{\eta}) = \mathbf{D}_\boldsymbol{\eta} A(\boldsymbol{\eta})^\top.$$

1.5.2.3 Bernoulli distribution

Definition 6. A scalar random variable x is from the Bernoulli distribution with probability of success p when its probability mass function is

$$f(x) = p^x(1-p)^{1-x}, \quad x \in \{0, 1\}, \quad (1.8)$$

which is written as $x \sim \text{Bernoulli}(p)$.

The sufficient statistic and base measure for the Bernoulli distribution are respectively defined as

$$T(x) = x \quad \text{and} \quad h(x) = I(x \in \{0, 1\}).$$

Definition 7. For a scalar input variable x

$$\text{logit}(x) \equiv \log\left(\frac{x}{1-x}\right) \quad \text{and} \quad \text{expit}(x) \equiv \frac{\exp(x)}{1+\exp(x)}. \quad (1.9)$$

The natural parameter vector and inverse mapping are respectively

$$\eta = \text{logit}(p) \quad \text{and} \quad p = \text{expit}(\eta). \quad (1.10)$$

The log-partition function is

$$A(\eta) = \log(1 + \exp(\eta)).$$

The domain of $A(\eta)$ and $\nabla A(\eta)$ is $H = \eta : \eta \in \mathbb{R}$. The derivative of $A(\eta)$ is

$$(\nabla A)(\eta) = \frac{\exp(\eta)}{1 + \exp(\eta)}.$$

1.5.2.4 Poisson distribution

Definition 8. A scalar random variable $x \in \mathbb{Z}_{\geq 0}$ is Poisson distributed with mean and variance λ when its density function is

$$f(x) = \frac{\lambda^x \exp(-\lambda)}{\Gamma(x+1)}, \quad (1.11)$$

which is written as $x \sim \text{Poisson}(\lambda)$, where $\Gamma(x) = (x-1)!$.

The sufficient statistic and base measure for the Poisson distribution are respectively defined as

$$T(x) = x \quad \text{and} \quad h(x) = I(x \in \mathbb{Z}_{\geq 0})$$

The natural parameter vector and inverse mapping are respectively

$$\eta \equiv \log \lambda \quad \text{and} \quad \lambda = \exp(\eta). \quad (1.12)$$

The log-partition function is

$$A(\eta) = \exp(\eta).$$

The domain of both $A(\eta)$ and $\nabla A(\eta)$ is $H = \eta : \eta \in \mathbb{R}$. The derivative of $A(\eta)$ is

$$(\nabla A)(\eta) = \exp(\eta).$$

1.5.2.5 Negative binomial distribution

Definition 9. A distributed scalar random variable $x \in \mathbb{Z}_{\geq 0}$ is from the negative binomial distribution with probability of success p and shape parameter $\kappa > 0$ when its density function is

$$f(x) = \frac{\Gamma(x + \kappa)}{\Gamma(x + 1)\Gamma(\kappa)} \left(\frac{\mu}{\mu + \kappa}\right)^x \left(\frac{\kappa}{\mu + \kappa}\right)^\kappa, \quad (1.13)$$

which is written as $x \sim NB(\mu, \kappa)$, where $\Gamma(x) = (x - 1)!$ and $p = \mu/(\mu + \kappa)$.

The negative binomial distributions is closely related to the Poisson distribution, and begins to resemble it as the shape parameter κ increases to infinity. The sufficient statistic and base measure for the multivariate normal family are respectively defined as

$$T(x) = x \quad \text{and} \quad h(x) = \binom{x + \kappa - 1}{x}.$$

The natural parameter vector and inverse mapping are respectively

$$\eta \equiv \log \mu \quad \text{and} \quad \mu = \exp(\eta). \quad (1.14)$$

The log-partition function is

$$A(\eta) = -\kappa \log(1 - \exp(\eta)).$$

The domain of both $A(\eta)$ and $\nabla A(\eta)$ is $H = \eta : \eta \in \mathbb{R}$. The derivative of $A(\eta)$ is

$$(\nabla A)(\eta) = \frac{\kappa \exp(\eta)}{1 - \exp(\eta)}.$$

1.6 Graph theory

Graph theory is vital in understanding message passing and helps to simplify what is an otherwise complex task. In this section we provide a brief review of the required graphical theory.

1.6.1 Directed acyclic graphs

Directed acyclic graphs are a graphical representation of a model consisting of nodes and edges. The nodes are used to represent vectors/matrices, while the lines between them known as edges, demonstrates the relationships between nodes. To be a directed acyclic graph, all edges must be directed and without cycles; that is, there cannot be any connected directed edge, where following the direction of said edge returns to the starting edge (see Figure 1.1b).

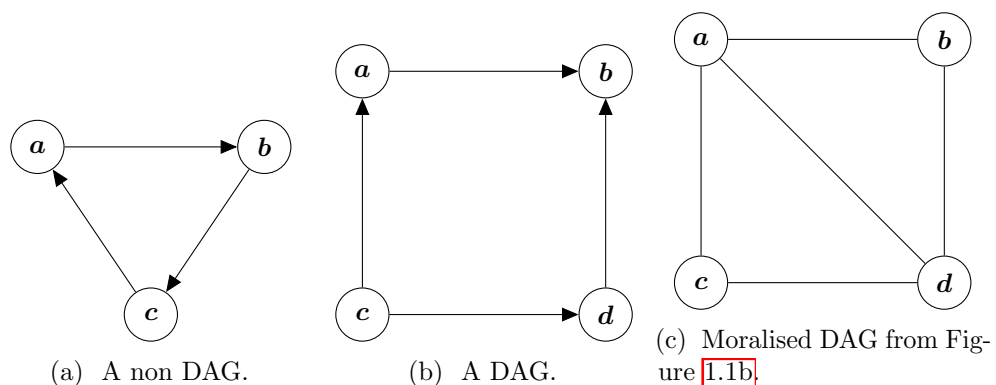


Figure 1.1: *Three basic graphical models. Figure 1.1b shows a directed acyclic graph (DAG); all edges are directed and without cycles. Conversely, Figure 1.1a is not a directed acyclic graph. Although it contains directed edges, it also contains a cycle. The moralised version of the directed acyclic graphs from Figure 1.1b is Figure 1.1c. No edges have a direction and the parent nodes a and d are linked.*

Moralisation of directed acyclic graphs help uncover the relationship between nodes. A parent node is the node connected away from the directed edge, while the node on the directed edge is the child node. Co-parents occur when a child node has two edges directed towards it. Specifically, the process of moralisation involves linking parents and co-parents. Consider the moralisation of Figure 1.1b: the child node b has co-parents a and d . To moralise this graph, we link nodes a and d with an undirected edge. However, no further links are required for the parent node of a and d (node c) since it does not have a co-parent. We then change all directed edges to undirected edges (Figure 1.1c).

The a node's Markov blanket consists of the nodes directly linked to it. Finding the Markov blanket is easy on a moralised graph and can lead to significant simplifications for large graphical models. For example, in our moralised graph, the Markov blanket of node a would be nodes b , c and d , where as the Markov blanket for node c would be nodes a and d (Pearl, 1988⁵³).

1.6.2 Factor graphs

An extension of moralised graphs are factor graphs. Figure 1.2 demonstrates a factor graph corresponding to the model in equation (1.15),

$$h(x_1, x_2, x_3, x_4, x_5) = \sqrt{x_1 + x_2 x_5} \log(x_2 + x_3^2 x_4^7)^{1/2} \sin(x_3) (|x_5|^2 - 10), \quad (1.15)$$

where the square nodes are functions in the equation known as factor nodes, and the circle nodes are arguments in the equation known as stochastic nodes. Neighbours of a

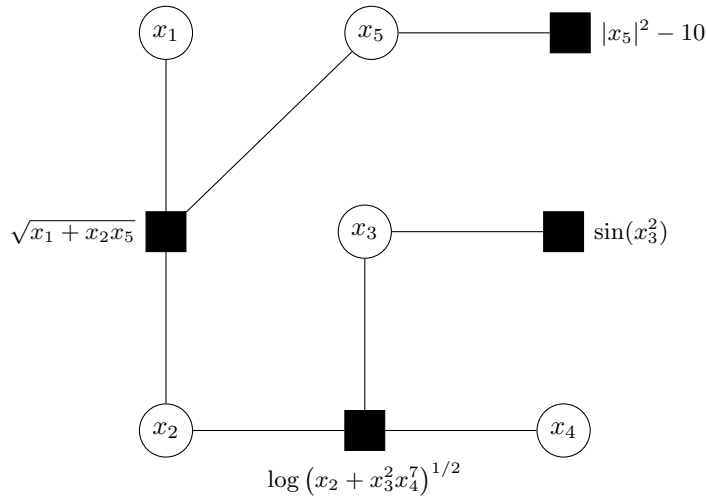


Figure 1.2: The factor graph corresponding to the model presented in equation (1.15). The two types of nodes in the factor graph are shown here, where square nodes represent factor nodes and circular nodes represent stochastic nodes.

node are the those nodes linked by an edge. For example, the factor node $\sqrt{x_1 + x_2 x_5}$ in Figure 1.2 is neighbours with stochastic nodes (x_1, x_2, x_5) , whereas the factor node $|x_5|^2 - 10$ only is neighbours with the stochastic node (x_5) (Rohde & Wand, 2015^[58]).

1.7 Multilevel datasets

Often it is not possible or practical to collect data from a single experimental group with one source of variance. As such, real world datasets are often structured in a multilevel or hierarchical manner, whereby they have nested levels stemming from observations occurring within groups (Steenbergen & Jones, 2002^[62]). For categorical variable A and B , A is nested in B if each category of A only occurs in one category of B . A variable is defined as a level if its values are a random sample from a wider population of values,

in which case the values it takes are called groups. In other words, groups of a level are a sample from a larger population of groups, where the observations have a distribution (McCulloch, Searle & Neuhaus, 2008^[41]). Alternatively, a variable is a level if there is a functional difference between its values (Borenstein, et al., 2009^[10]). Each level of the structure contains an unobserved residual component, which accounts for correlation within the structure.

Levels are a consequence of experimental design and data collection. We aim to clarify what constitutes a level by explaining an example. Before we continue with the example, we note variables can have two types of effect on the response. Fixed effects occur when a variable is deemed to have a constant functional effect on the response from a finite set of functional effects. Generally analysts are aware of such variables in studies. Random effects occur when variables are deemed to have differing functional effects from an infinite set of functional effects. They are assumed to be a sample of a random categorical variables from some distribution and are caused by unknown variables (McCulloch, Searle & Neuhaus, 2008^[41]). It is often not immediately obvious when variables have fixed or random effects, with inference being a key motivating factor. Additionally there is some overlap between variables that are levels and cause random effects. We aim to clarify this in our example as well.

Consider a school that wants to compare the effect of time spent studying on pupils' exam scores for different subjects. Because class sizes are limited the school must collect pupils' exam scores from four classes (Raudenbush, 1993^[57]). The structure of this data is called a one level dataset and is shown in Figure [1.3](#), where pupils are nested in classes. Time spent studying is a continuous variable and no assumption is made that the number of hours studied are random, so it is likely similar values will have a similar effect on exam score. With this in mind the effect of time spent studying on pupils' exam scores is fixed and it is not level in the dataset. Although the effect of subject on exam score is not continuous or random, even in the case where subjects are chosen at random, it is more than likely their effects are explainable and should be considered fixed effects. This is also evidence they are not a level in the study.

As mentioned earlier, inference is an important consideration when determining levels, as well as fixed and random effects. If inference from the study is restricted to the four classes in the study only, the classes represent the population of classes in the study. In other words, they have a finite set of fixed functional effects and are not a level in the dataset (they should be included as a fixed effect). Alternatively, if inference from the study is extended to the population of students in the school, it is reasonable to assume classes in the study are a random sample from the wider population of classes in the

school (McCulloch, Searle & Neuhaus, 2008^[41]). These classes come from a probability distribution regarding their relationship with exam score. As such, in this case class is a level with four groups (one for each class). It should be clear that experimental design, data interpretation and inference all play a role in determining whether variables are a level and the type of effect they have on the response.

In this thesis we cover two other types of multilevel datasets which we now explain briefly with examples. The natural extension of the one level dataset explained above is a two level dataset, where an additional level of nesting is added. Continuing the exam score example, consider exam scores measured on pupils of different classes from multiple schools in an area. Assuming inference is now extended to the population of pupils in all classes of schools in the area, it is reasonable to assume both schools and classes have random effects with their own distributions. This is a two level dataset where pupils are nested in classes nested in schools, as shown in Figure 1.4.

Alternatively, datasets with two levels can be crossed. For two arbitrary levels B and C are crossed if each observation nested from each groups in level C occur in different groups of level B . Following our exam score example, consider exam scores are measured on pupils of different classes in the same school, however the pupils of each class live in different areas. Assuming inference is extended to both classes and the areas where pupils live, a crossed dataset results as shown in Figure 1.5, where pupils are nested in classes and classes crossed with area. As before, both levels have random effects with their own distributions. Note, although it the crossed example has two levels, the first level is not nested in the second level. Instead, they both nest pupils and can be thought of as structurally at the same level in the data.

Because multilevel datasets are composed of non-homogeneous experimental units we must model them accordingly. Failing to do so can facilitate incorrect model interpretation. Linear mixed models provide a solution to modeling multilevel datasets composed of non-homogeneous experimental units. They are a form of regression model which gained popularity with the rise of computing power. These models improve on linear models by allowing an intercept and slope for the population effects (known as fixed effects) and a unique intercept and slope to account for the effect of each experimental unit (known as random effects).

Perhaps the clearest advantage of mixed models is when they are used for prediction (Gelman, 2006^[18]). Mixed models allow for the estimation of group effects simultaneously with the effects of group-level predictors. This is not possible in fixed effects and ANOVA models, where the inferences are limited to the groups in the sample. With

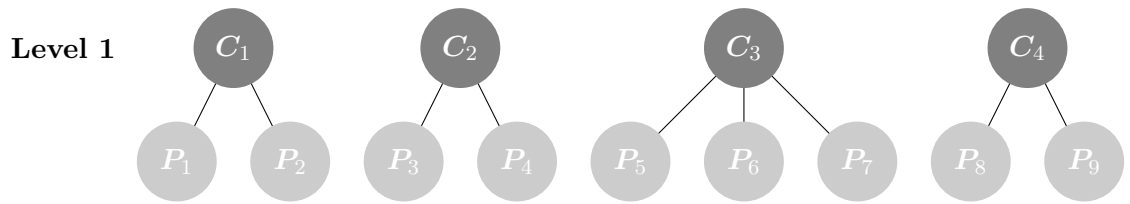


Figure 1.3: *Example of a two level dataset following the schools example. We denote P_j as the pupils in the classes and C_i the classes in the schools.*

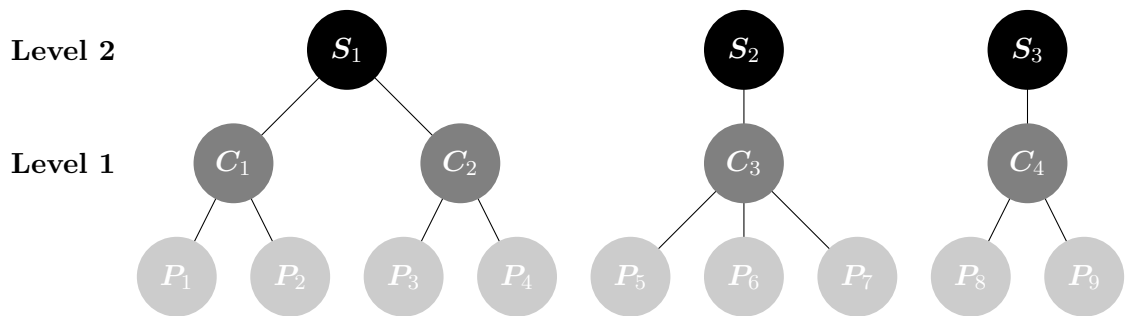


Figure 1.4: *Example of a two level dataset following the schools example. We denote P_j as the pupils in the classes, C_i the classes in the schools and S_k the schools.*

these traits in mind, it is clear that multilevel models provide a significant advantage over traditional regression techniques for prediction and data reduction.

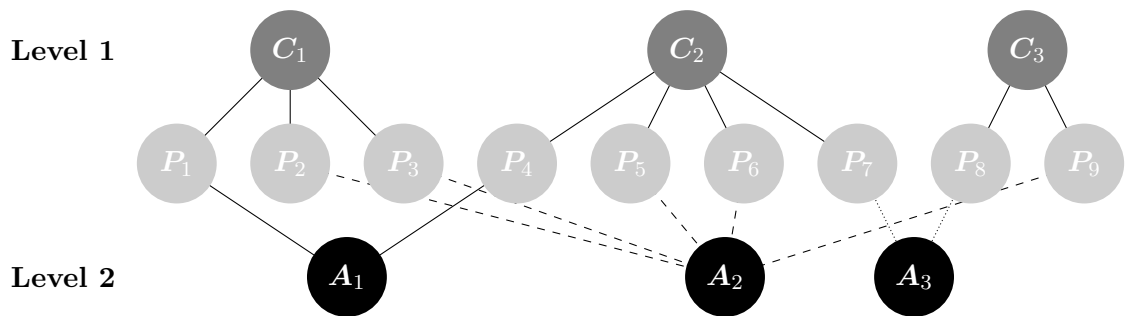


Figure 1.5: *Example of a crossed dataset following the schools example. We denote P_j as the pupils in the classes, C_i the classes in the schools and A_k the area.*

1.8 Generalised linear mixed models

Generalised linear mixed models are a powerful model that facilitate inference from multilevel datasets with non-normal response variables. Much like generalised linear models, a linear predictor η is used to incorporate information about the independent variables into the model. However, GLMMs include fixed and random effects in the linear predictor to account for error structures in the data, similar to linear mixed models. For GLMMs it is optional whether fixed intercepts and slopes are included in the linear predictor, however they must always include a random effect, via either an intercept or slope. The relationship between these model types is more clearly demonstrated in Table 1.1. The expected value of the response y conditional on the random effects is related to a linear predictor using a link function $g(\cdot)$, where $F(\cdot) = g^{-1}(\cdot)$ is the inverse link function, also called the mean function. In other words, although the mean is not directly a linear combination of predictors, some function of the mean is. We now show general models for the three data structures explained in the previous subsection.

For a dataset with one level of nesting analogous to Figure 1.3, where level one has

Table 1.1: *The relationship between four types of linear models, where LM refers to linear models, GLM refers to generalised linear models, LMM refers to linear mixed models and GLMM refers to generalised linear mixed models.*

	Response is normally distributed	Response is not normally distributed
Fixed effects only	LM	GLM
Fixed and random effects	LMM	GLMM

m groups and n_i measurements in each group, assuming the effect of each group is normally distributed and the distribution of the response y_{ij} conditional on the random effects is from an exponential family f , the mixed model has the general form

$$y_{ij}|\mathbf{u}_i \stackrel{\text{ind.}}{\sim} f(F(y_{ij}|\mathbf{u}_i)), \quad \mathbf{u}_i \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R}}}, \boldsymbol{\Sigma}),$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i. \quad (1.16)$$

The linear predictor is related to the expected value of the response by the link function F

$$g(E(y_{ij}|\mathbf{u}_i)) = \eta_{ij}, \quad \text{where} \quad E(y_{ij}|\mathbf{u}_i) = F(\eta_{ij}).$$

For the one level mixed model the linear predictor is

$$\eta_{ij} = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}},$$

where $\mathbf{x}_{ij}^{\mathbf{F}}$ is a $d^{\mathbf{F}} \times 1$ vector of predictors modelled as having fixed effects with coefficient vector $\boldsymbol{\beta}$, and $\mathbf{x}_{ij}^{\mathbf{R}}$ is a $d^{\mathbf{R}} \times 1$ vector of predictors modelled as having random effects with coefficient vectors \mathbf{u}_i . Note, when the first entry of $\mathbf{x}_{ij}^{\mathbf{F}}$ or $\mathbf{x}_{ij}^{\mathbf{R}}$ equals 1, it corresponds to including a fixed or random intercept respectively.

For a dataset with two levels of nesting analogous to Figure [1.4](#), consisting of m outer groups, n_i inner groups in each outer group, and o_{ij} measurements each in each inner group, assuming the random effects $\mathbf{u}_i^{\mathbf{L}1}$ of the inner group and $\mathbf{u}_{ij}^{\mathbf{L}2}$ of the outer group are independent and identically distributed as normal, and the distribution of the response y_{ijk} conditional on the random effects is from an exponential family f , the mixed model has the general form

$$y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2} \stackrel{\text{ind.}}{\sim} f(F(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})),$$

$$\mathbf{u}_i^{\mathbf{L}1} \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R}1}}, \boldsymbol{\Sigma}^{\mathbf{L}1}) \quad \text{independently of} \quad \mathbf{u}_{ij}^{\mathbf{L}2} \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R}2}}, \boldsymbol{\Sigma}^{\mathbf{L}2}),$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij}. \quad (1.17)$$

The linear predictor is related to the expected value of the response by the link function F

$$g(E(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})) = \eta_{ijk} \quad \text{where} \quad E(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) = F(\eta_{ijk}).$$

For the two level mixed model the linear predictor is

$$\eta_{ijk} = \boldsymbol{\beta}^\top \mathbf{x}_{ijk}^{\mathbf{F}} + (\mathbf{u}_i^{\mathbf{L}1})^\top \mathbf{x}_{ijk}^{\mathbf{R}1} + (\mathbf{u}_{ij}^{\mathbf{L}2})^\top \mathbf{x}_{ijk}^{\mathbf{R}2},$$

where $\mathbf{x}_{ijk}^{\mathbf{F}}$ is a $d^{\mathbf{F}} \times 1$ vector of predictors, modelled as having fixed effects with coefficient vector $\boldsymbol{\beta}$, $\mathbf{x}_{ijk}^{\mathbf{R}1}$ is a $d^{\mathbf{R}1} \times 1$ vector of predictors modelled as having random effects from

the outer groups with coefficient vectors $\mathbf{u}_i^{\mathbf{L}1}$, and $\mathbf{x}_{ijk}^{\mathbf{R}2}$ is a $d^{\mathbf{R}2} \times 1$ vector of predictors modelled as having random effects from the inner groups with coefficient vectors $\mathbf{u}_{ij}^{\mathbf{L}2}$. Note, when the first entry of $\mathbf{x}_{ijk}^{\mathbf{F}}$, $\mathbf{x}_{ijk}^{\mathbf{R}1}$ or $\mathbf{x}_{ijk}^{\mathbf{R}2}$ equals 1, it corresponds respectively to including a fixed or random intercept for the outer or inner groups.

For a dataset with crossed levels of nesting analogous to Figure 1.5, consisting of two levels of m and m' groups, with observations indexed according to the pair $(i, i') \in \{1, \dots, m\} \times \{1, \dots, m'\}$, and $n_{ii'}$ observations in the (i, i') pair, assuming the random effects of each level denoted by \mathbf{u}_i and $\mathbf{u}'_{i'}$ are independent and identically distributed as normal, and the distribution of the response $y_{ii'j}$ conditional on the random effects is from an exponential family f , the crossed mixed model has the general form

$$\begin{aligned} y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'} &\stackrel{\text{ind.}}{\sim} f(F(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'})), \\ \mathbf{u}_i &\stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R}}}, \boldsymbol{\Sigma}) \quad \text{independently of} \quad \mathbf{u}'_{i'} \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R}'}} , \boldsymbol{\Sigma}'), \\ 1 \leq i \leq m, \quad 1 \leq i' \leq m', \quad 1 \leq j \leq n_{ii'}. \end{aligned} \quad (1.18)$$

The linear predictor is related to the expected value of the response by the link function F

$$g(E(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'})) = \eta_{ii'j} \quad \text{where} \quad E(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}) = F(\eta_{ii'j}).$$

For the crossed mixed model the linear predictor is

$$\eta_{ii'j} = \boldsymbol{\beta}^\top \mathbf{x}_{ii'j}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ii'j}^{\mathbf{R}} + (\mathbf{u}'_{i'})^\top \mathbf{x}_{ii'j}^{\mathbf{R}'},$$

where $\mathbf{x}_{ii'j}^{\mathbf{F}}$ is a $d^{\mathbf{F}} \times 1$ vector of predictors, modelled as having fixed effects with coefficient vector $\boldsymbol{\beta}$, $\mathbf{x}_{ii'j}^{\mathbf{R}}$ is a $d^{\mathbf{R}} \times 1$ vector of predictors modelled as having random effects from the groups of the first level with coefficient vectors \mathbf{u}_i , and $\mathbf{x}_{ii'j}^{\mathbf{R}'}$ is a $d^{\mathbf{R}'} \times 1$ vector of predictors modelled as having random effects from the groups of the second level with coefficient vectors $\mathbf{u}'_{i'}$. Note, when the first entry of $\mathbf{x}_{ii'j}^{\mathbf{F}}$, $\mathbf{x}_{ii'j}^{\mathbf{R}}$ or $\mathbf{x}_{ii'j}^{\mathbf{R}'}$ equals 1, it corresponds respectively to including a fixed intercept, or a random intercept for the first or second groups.

Although other choices exist, selecting the link function F as the canonical link for the conditional response distribution f is a good default and provides a number of desirable properties. We now present distributions and link functions for binary and count response data on an arbitrary model structure. The distributions and links that follow can be applied to the random effects structures presented.

1.8.1 Binary response models

When dealing with binary response variables where $y \in \{0, 1\}$, it is common to assume the response is Bernoulli distributed where

$$y|u \sim \text{Bernoulli}(F(\eta)).$$

For binary outcomes, several popular inverse link functions exist, however we focus on two:

$$F = \begin{cases} \text{expit} & \text{logistic inverse link,} \\ \Phi & \text{probit inverse link} \end{cases}$$

where expit and Φ are defined in Section [1.5.2](#). The expit inverse link (inverse of the logit link) function is the most frequently used link when modelling binary data. It affords an elegant log-odds model interpretation as well as algebraic simplifications due to it being the canonical link of the Bernoulli distribution. Although Probit links often facilitate more tractable solutions, their model coefficients do not have a direct interpretation in the same way that logit links do, and inference is restricted. As a side note, although special cases of grouped data with responses proportional to the number of observations per group exist, they are uncommon and as such are not discussed in this thesis.

1.8.2 Count response models

When modelling count variables where $y \in \mathbb{Z}_{\geq 0}$, two models are commonly used. The Poisson model assumes the response is from a Poisson distribution as

$$y|u \stackrel{\text{ind.}}{\sim} \text{Poisson}(F(\eta)),$$

where the inverse link is the canonical function $F = \exp$. As Poisson distributed variables are constrained to have an equal mean and variance, it is not appropriate for over-dispersed data where the variance is greater than the mean.

The negative binomial model is similar to the Poisson model but includes an additional parameter κ to account for over-dispersed data. It can be written as

$$y|u \stackrel{\text{ind.}}{\sim} \text{NB}(F(\eta), \kappa),$$

where as before the inverse link is the canonical function $F = \exp$. As the shape parameter $\kappa \rightarrow \infty$ the negative binomial distribution converges to the Poisson distribution.

1.9 Maximum likelihood

1.9.1 Likelihood functions

The likelihood of a model is the probability of observing the data given parameter inputs. It is a function of the parameters where random variables are fixed at the observed values. For a statistical model with probability density function $f(y)$ for data y and unknown parameters $\boldsymbol{\theta}$, the likelihood function is equivalent to the probability density function of the observed data y ,

$$L(\boldsymbol{\theta}; y) = f(y; \boldsymbol{\theta}),$$

with a log-likelihood function given by (Collins, 2008^[16])

$$\ell(\boldsymbol{\theta}; y) = \log L(\boldsymbol{\theta}).$$

In the case of m independent observations of y_i

$$\ell(\boldsymbol{\theta}; y) = \sum_{i=1}^m \log f(y_i; \boldsymbol{\theta}).$$

Note that by using the log-likelihood we sum over the independent observations as opposed to taking their product, which helps reduce numerical instability.

For the general one level model given in equation (1.16), the joint log-likelihood function for parameters $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sum_{i=1}^m \log \int_{\mathbb{R}^{d\mathbf{R}}} \prod_{j=1}^{n_i} p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) p(\mathbf{u}_i; \boldsymbol{\Sigma}) d\mathbf{u}_i,$$

where

$$p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) = f(F(y_{ij}|\mathbf{u}_i)) \quad \text{and} \quad p(\mathbf{u}_i; \boldsymbol{\Sigma}) = N(\mathbf{0}_{d^{\mathbf{R}}}, \boldsymbol{\Sigma}).$$

For the general two level model given in equation (1.17), the log-likelihood for the parameters $(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}})$ is

$$\begin{aligned} \ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}) &= \sum_{i=1}^m \log \int_{\mathbb{R}^{d^{\mathbf{R1}}+d^{\mathbf{R2}}}} \prod_{j=1}^{n_i} \prod_{k=1}^{o_{ij}} p \left(y_{ijk} \mid \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}; \boldsymbol{\beta} \right) \\ &\quad \times p \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}; \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}} \right) d \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}, \end{aligned}$$

where

$$p \left(y_{ijk} \mid \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}; \boldsymbol{\beta} \right) \equiv f(F(y_{ijk}|\mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}}))$$

and

$$p \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}; \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}} \right) \equiv N \left(\begin{bmatrix} \mathbf{0}_{d^{\mathbf{R1}}} \\ \mathbf{0}_{d^{\mathbf{R2}}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}^{\mathbf{L1}} & \mathbf{0}_{d^{\mathbf{R1}} \times d^{\mathbf{R2}}} \\ \mathbf{0}_{d^{\mathbf{R2}} \times d^{\mathbf{R1}}} & \boldsymbol{\Sigma}^{\mathbf{L2}} \end{bmatrix} \right).$$

For the general crossed random effects model given in equation (1.18), the log-likelihood for the parameters $(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ may be written as

$$\begin{aligned} \ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') &= \log \int_{\mathbb{R}^{m d^{\mathbf{R}}+m' d^{\mathbf{R}'}}} \prod_{(i,i'):n_{ii'}>0} \prod_{j=1}^{n_{ii'}} p \left(y_{ii'j} \mid \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}; \boldsymbol{\beta} \right) \\ &\quad \times p \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}' \right) d \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}, \end{aligned}$$

where

$$p \left(y_{ii'j} \mid \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}; \boldsymbol{\beta} \right) \equiv f(F(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}))$$

and

$$p \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}' \right) \equiv N \left(\begin{bmatrix} \mathbf{0}_{m d^{\mathbf{R}}} \\ \mathbf{0}_{m' d^{\mathbf{R}'}} \end{bmatrix}, \begin{bmatrix} \mathbf{I}_m \otimes \boldsymbol{\Sigma} & \mathbf{0}_{m d^{\mathbf{R}} \times m' d^{\mathbf{R}'}} \\ \mathbf{0}_{m' d^{\mathbf{R}'}} \times m d^{\mathbf{R}} & \mathbf{I}_{m'} \otimes \boldsymbol{\Sigma}' \end{bmatrix} \right).$$

1.9.2 Maximum likelihood

Parameters at the peak of a likelihood function hypersurface can be interpreted as those that give the maximum probability of observing the data obtained under the fitted model. This process forms a frequentist probabilistic framework called maximum likelihood estimation. By taking $L(\boldsymbol{\theta}|y)$ as a function of $\boldsymbol{\theta}$, the goal of maximum likelihood estimation is to return the parameters that maximise the $L(\boldsymbol{\theta}|y)$ (denoted by $\hat{\boldsymbol{\theta}}$).

Generally, maximum likelihood estimates can be solved by differentiating the likelihood function with respect to parameters and solving to find the root. The maxima of the likelihood and log-likelihood are equivalent, although as mentioned it is often preferable to work with the log-likelihood for several reasons. As such, the maximum likelihood estimates are given by

$$\left. \frac{d\ell}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = \mathbf{0}_d.$$

Maximum likelihood estimates are consistent and asymptotically normally distributed, with mean equal to $\boldsymbol{\theta}$ and variance equal to the inverse Fisher information (McCulloch, Searle & Neuhaus, 2008^[41]),

$$\hat{\boldsymbol{\theta}} \sim \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}(\hat{\boldsymbol{\theta}})^{-1}).$$

Two important issues arise from this. Firstly, the range of parameters being estimated must be considered. Specifically, in the case of variance parameters which occur only on the positive real number line, the search space must be constrained. Second, analytic solutions are not always tractable or practical to obtain. For models with parameters of varying dimensions derivative free options become more appealing.

1.10 Best prediction

GLMMs are often used for predicting values of random effects, despite the fact they are unobservable. The best predictor (BP) is the predictor with the lowest mean squared error. That is, for an arbitrary GLMM with data y_{ij} , the best predictor of the random effect \mathbf{u}_i is

$$\text{BP}(\mathbf{u}_i) \equiv \underset{\mathbf{u}_i \in \mathbb{R}}{\text{argmin}} E[(\mathbf{u}_i - \mathbf{u})^2].$$

It is straightforward to show the best predictor is the average of the random effect (Scott, Simonoff & Marx, 2013^[60])

$$\text{BP}(\mathbf{u}_i) = E(\mathbf{u}_i | y_{ij}).$$

This can be calculated using the conditional specification

$$\text{BP}(\mathbf{u}_i) = \int_{\mathbb{R}^{d_{\mathbf{R}}}} \mathbf{u}_i \left(\frac{f(y_{ij} | \mathbf{u}_i) f(\mathbf{u}_i)}{\int_{\mathbb{R}^{d_{\mathbf{R}}}} f(\mathbf{y}_i | \mathbf{u}_i) f(\mathbf{u}_i)} \right) d\mathbf{u}_i.$$

McCulloch & Neuhaus (2012)^[42] show prediction accuracy measured by mean square error is robust to moderate violations of the assumed random effects distribution. This suggests that for prediction, inference is relatively immune to variations of assumptions. Since it is difficult to confirm whether or not the assumptions of the random effects distribution are met, this attribute of best predictors is extremely useful.

A number of prediction schemes such as best linear predictors and best linear unbiased predictors. Best linear predictors are useful when the full PDF is not available as it only requires the first and second order moments. Unlike best predictors, restrictive assumptions are imposed in best linear predictors to minimise the mean squared error. Although these assumptions may result computation of results, they also ensure that best linear predictors never has a smaller mean squared error than best predictors. Best linear unbiased predictors are a subclass of best linear predictors, and as such their mean squared error never has a smaller mean squared error than best linear predictors. Sections 2 to 3 of Teunissen (2007)^[63] provide a good overview of the differences between predictors and we refer interested readers there.

1.11 Current approximation methods

The integrals involved with calculating the likelihood surface and best predictors of GLMMs do not have tractable solutions. In the frequentist GLMM setting the two standard approaches to solving intractable integrals are Laplace approximations and Gauss-Hermite quadrature. We now provide a brief overview of each method.

1.11.1 Laplace approximation

Laplace approximations were established in 1774 and have since gone on to become a fundamental technique in mathematics and statistics. These approximations are in effect a second order Taylor series approximation around the maximum of the function to be approximated. This approximation allows the integral to be expressed as a Gaussian distribution with an analytic solution.

Consider an integrals of the form

$$I(x) = \int_a^b \exp(f(x)) dx,$$

where the function $f(x)$ has a maxima at x_0 such that $a < x_0 < b$ and $f''(x_0) < 0$. A second order Taylor expansion of $f(x)$ around x_0 is

$$f(x) \approx f(x_0) + (x - x_0)f'(x_0) + \frac{1}{2}(x - x_0)^2 f''(x_0).$$

Since $f'(x_0) = 0$

$$f(x) = f(x_0) - \frac{1}{2}(x - x_0)^2 (-f''(x_0)).$$

It can be seen that $\exp(f(x))$ is expressible as a normal density function with mean x_0 and variance $-f''(x_0)^{-1}$. Given the assumptions imposed on $f(x)$

$$\int_a^b \exp(f(x)) dx \approx \exp(f(x_0)) \int_a^b \exp\left(\frac{(x - x_0)^2}{2f''(x_0)^{-1}}\right) dx.$$

When the integral on the right side of the previous equation is evaluated between $-\infty$ and ∞ it becomes a Gaussian integral, thus

$$\int_a^b \exp(f(x)) dx \approx \exp(f(x_0)) \sqrt{-2\pi f''(x_0)^{-1}}.$$

For an arbitrary GLMM with log-likelihood function of parameters $\boldsymbol{\theta}$

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^m \ell_i(\boldsymbol{\theta})$$

where

$$\ell_i(\boldsymbol{\theta}) = \int_{\mathbb{R}^{d\mathbf{R}}} \prod_{j=1}^{n_i} f(y_{ij}|\mathbf{u}_i) f(\mathbf{u}_i) d\mathbf{u}_i,$$

the integral can be expressed as

$$\begin{aligned}\ell_i(\boldsymbol{\theta}) &= \int_{\mathbb{R}^d} \exp \left\{ \sum_{j=1}^{n_i} \log (f(y_{ij}|\mathbf{u}_i)f(\mathbf{u}_i)) \right\} d\mathbf{u}_i \\ &= \int_{\mathbb{R}^d} \exp \tilde{f}(\mathbf{u}_i) d\mathbf{u}_i,\end{aligned}$$

where

$$\tilde{f}(\mathbf{u}_i) = \sum_{j=1}^{n_i} \log f(y_{ij}|\mathbf{u}_i)f(\mathbf{u}_i).$$

It can subsequently be solved as shown above. Laplace approximation provides fast inference for GLMMs, performing well in cases where the number of observations of each group is large, i.e. large n_i . However, they are less accurate in cases where the number of observations of each group are low and the random effects have high variance. As such, it is worthwhile considering alternative approximation methods.

1.11.2 Gauss Hermite quadrature

Perhaps the simplest and most robust method of numerical integration, quadrature facilitates an approximate calculation of definite integrals. Gauss-Hermite quadrature (GHQ) approximates definite integrals from the normal family or with a log-quadratic factor by a weighted sum. For polynomials of degree $2m - 1$ GHQ provides an exact result or less over the domain $[-1, 1]$ by focusing on selection of m optimal nodes x_i . Each node corresponds to the roots of an m th order Hermite polynomial $H_m(x)$ and is accompanied by a weight w_i defined by the Gauss-Hermite weight function (Liu & Pierce, 1994^[36])

$$w_i = \frac{2^{m-1} m! \sqrt{\pi}}{m^2 (H_{m-1}(x_i))^2}.$$

The nodes are symmetric about zero, where the range increases with m . The integral is then computed as the weighted sum of function values $f(x_i)$ at these points

$$\int_{-\infty}^{\infty} f(x)\phi(x)dx = \sum_{i=1}^m w_i f(x_i).$$

Note, that since Gaussian quadrature works over the domain $[-1, 1]$, a change of interval is required from $[a, b]$. Implementation of the Gaussian quadrature rule suggests the

integral can be approximated as follows

$$\int_a^b f(x)\phi(x)dx \approx \frac{b-a}{2} \sum_{i=1}^n w_i f\left(\frac{b-a}{2} x_i + \frac{a+b}{2}\right).$$

Adaptive quadrature follows a similar process to the one for traditional quadrature, however also implements an accuracy criterion based on the difference between two subintervals. If a large discrepancy exists between two intervals (i.e. there is a large amount of error between them), the subinterval is re-divided into two smaller subintervals and recalculated. It is this process which improves the accuracy of adaptive quadrature over traditional methods for poorly behaved functions. Additionally, the computational burden of this refinement process is reduced by implementing the accuracy criterion only when needed for smaller subintervals. The R function “`integrate()`” in the “stats” package allows for easy implementation of the adaptive quadrature via the Gauss-Kronrod method. Gauss-Kronrod quadrature extends Gaussian quadrature by adding $n + 1$ points to an n -point rule, such that the resulting rule is a polynomial of degree $2n + 1$. As such, a set of function evaluation points can be created, where the extra points (corresponding to the Kronrod extension) allow the computation of higher-order estimates, and the function values at the set points (corresponding to the Gaussian quadrature rule) provide lower-order estimates. The difference between these values forms the accuracy criterion previously discussed.

A major advantage of GHQ over other methods of quadrature is that once weights and nodes are calculated they can be stored to reduce computational costs of future calculations. For GLMMs the complexity of GHQ increases with the dimension of random effects. It additionally becomes inaccurate if the dimension of the random effects is greater than two due to limitations in the method used to factorise the high-dimensional integrals into a lower-dimensional one. As such it is not suitable for crossed models and models with more than two levels (Handayani, et al., 2017^[27]).

1.11.3 Other methods

In addition to the previous a number of other methods exist. Penalised quasi-likelihood (PQL) (Breslow & Clayton, 1993^[11]) is one of the most popular and high speed methods for handling the difficult integrals arising in GLMMs. PQL is useful in cases where we are missing information regarding the distribution of data and as such cannot obtain a full likelihood function, in which case a quasi-likelihood function can be used with a

penalty on the random effects component. The penalty ensures values of the random effects have a mean of zero. Laplace approximation can then be used to solve the intractable integral which arises. Although it provides fast inference and is robust to model misspecification it is not accurate and provides biased estimates when there is little data per group, such as in binary or low count data cases. More importantly, PQL calculates a quasi-likelihood rather than the true likelihood and as such may not be appropriate in situations where likelihood ratio testing is implemented.

Although traditional methods such as Markov chain Monte Carlo provide highly reliable estimates of intractable integrals for Bayesian models, there is little literature regarding implementation in a frequentist framework. Lele, et al. (2007)^[55] present a reformulation of Markov chain Monte Carlo, known as data cloning, which allows maximum likelihood and confidence interval calculation. This method involves building a fully specified Bayesian model of the problem with uninformative priors and creating a large number of copies of the data which are assumed to be independent. The posterior is then calculated with the usual Markov chain Monte Carlo approach and the likelihood over the copies is used as the data. The mean of the resulting posterior distribution is equal to the maximum likelihood estimate and the number of copies times the variance of the posterior is equal to the variance of maximum likelihood estimate. Unlike GHQ and Laplace approximation, it easily extends to consider multiple random effects. However, this approach involves several difficult technical details. Additionally, although heuristic framework for data cloning exist, current software packages do not yet support the looping facility required for its fitting. Furthermore, it is not widely proven and still bares inherent issues of Markov chain Monte Carlo algorithms such as being computationally intensive.

1.12 Expectation propagation

Expectation propagation (EP) is a Bayesian iterative algorithm used for the computation of intractable posterior distributions. While Opper & Winther (2000)^[52] first provided the scheme for Gaussian approximating families, Minka (2001)^[43] provided the generalised scheme of the EP algorithm for all exponential families. EP is a reinterpretation of assumed density filtering (Opper, 1999^[51]) such that the posterior approximation is no longer dependent on the order of data. Although this makes EP

more computationally intensive than assumed density filtering, EP is more accurate and consistent. Both algorithms are examples of mean field approximations (a common type of variational inference). These algorithms rely on approximating intractable posteriors by selecting approximate densities from tractable families by minimising a distance measure. This divergence measure and the choice of approximating density family differentiate methods of mean field approximations. The Kullback Leibler (KL) divergence is a popular distance measure used in mean field approximations, which we now explain.

1.12.1 Kullback Leiber divergence and projection

The KL divergence is defined as a measure of distance between two density functions f_1 and f_2 such that

$$\text{KL}(f_1\|f_2) = \int_{\mathbb{R}^d} f_1(\mathbf{x}) \log \left(\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right) + f_2(\mathbf{x}) - f_1(\mathbf{x}) d\mathbf{x}, \quad (1.19)$$

where both densities are on \mathbb{R}^d . In other words, a KL divergence of zero indicates the two density functions are identical. The two rightmost terms of equation (1.19) are a correction factor for unnormalised densities, thus in the case of normalised densities the KL divergence can be simplified to

$$\text{KL}(f_1\|f_2) \equiv \int_{\mathbb{R}^d} f_1(\mathbf{x}) \log \left(\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right) d\mathbf{x}.$$

By Gibb's inequality it can be shown that $\text{KL}(f_1\|f_2) \geq 0$. It is also easy to show that KL divergence is an asymmetric measure where $\text{KL}(f_1\|f_2) \neq \text{KL}(f_2\|f_1)$.

The KL projection of the density function f onto a family Q is defined as the distribution q closest to f in the family of density functions Q ,

$$\text{proj}[f] \equiv \underset{q \in Q}{\text{argmin}} \text{KL}(f\|q). \quad (1.20)$$

This is a difficult problem to solve without placing any constraints on the family Q . However, by constraining Q to be an exponential family density function the problem is

simplified to moment-matching. Suppose

$$Q = \left\{ q : q(\mathbf{x}) = \exp(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}))h(\mathbf{x}), \quad \boldsymbol{\eta} \in \mathbf{H} \right\}.$$

Then projection of f onto q is solved by

$$\text{proj}[f] = \exp(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta}^* - A(\boldsymbol{\eta}^*))h(\mathbf{x}),$$

where $\boldsymbol{\eta}^*$ is

$$\boldsymbol{\eta}^* = \underset{\boldsymbol{\eta} \in \mathbf{H}}{\text{argmin}} \left(A(\boldsymbol{\eta}) - \boldsymbol{\eta}^\top \int_{-\infty}^{\infty} \mathbf{T}(x)f(x)dx \right).$$

However, since we are working with the exponential family, it can be shown that the derivative vector of $A(\boldsymbol{\eta})$ equates to the expectation of the natural statistic and thus $\boldsymbol{\eta}^*$ is the solution to

$$\int_{-\infty}^{\infty} \mathbf{T}(x)f(x)dx = \int_{-\infty}^{\infty} \mathbf{T}(x) \exp(\mathbf{T}(x)^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}))h(x)dx.$$

That is, $\boldsymbol{\eta}^*$ is chosen such that f and $\text{proj}[f]$ have the same natural statistic moments. When the approximating family is Gaussian, Result 1 of Kim & Wand (2016)^[31] follows:

Result 1. *Let x be a non-degenerate random variable with density function f . The KL projection of a density function f onto the normal family, denoted by $\text{proj}_N[f]$, is the $N(\boldsymbol{\mu}^*, (\boldsymbol{\sigma}^2)^*)$ density function where*

$$\boldsymbol{\mu}^* = E(x) \quad \text{and} \quad (\boldsymbol{\sigma}^2)^* = E(x^2) - (E(x))^2.$$

Result [1] is easily extended to the multivariate normal distribution as in Result [2]. That is, the projection is chosen to be the d -variate normal density function with the same mean vector and covariance as f (Kim & Wand, 2017^[68]).

Result 2. *Let \mathbf{x} be a non-degenerate random variable with density function f . The KL projection of a d -variate density function f onto the multivariate normal family, denoted by $\text{proj}_N[f]$, is the $N(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$ density function where*

$$\boldsymbol{\mu}^* = E(\mathbf{x}) \quad \text{and} \quad \boldsymbol{\Sigma}^* = E(\mathbf{x}\mathbf{x}^\top) - E(\mathbf{x})E(\mathbf{x})^\top.$$

1.12.2 Mean field approximations

We now provide a brief overview of mean field approximations and how popular variational methods such as mean field variational Bayes and belief propagation are related to EP. For observed data \mathbf{x} and parameter vector $\boldsymbol{\theta}$ consider approximations to the joint posterior density function $f(\boldsymbol{\theta}|\mathbf{x})$ of the form

$$f(\boldsymbol{\theta}|\mathbf{x}) \approx q^*(\boldsymbol{\theta}),$$

where $q^*(\boldsymbol{\theta}) = \prod_{i=1}^m q^*(\boldsymbol{\theta}_i)$ and $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m)$ is a partition of $\boldsymbol{\theta}$ and the $q^*(\boldsymbol{\theta}_i)$ s are chosen to minimise the KL divergence of $f(\boldsymbol{\theta}|\mathbf{x})$ from a product density over the elements of the partition, that is

$$q^*(\boldsymbol{\theta}) = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \left(\prod_{i=1}^m q(\boldsymbol{\theta}_i) \parallel f(\boldsymbol{\theta}|\mathbf{x}) \right). \quad (1.21)$$

The optimisation problem in equation (1.21) corresponds to mean field variational Bayes. This leads to a mode seeking approximation, where $q^*(\boldsymbol{\theta})$ is selected by maximising the probability of being under $f(\boldsymbol{\theta}|\mathbf{x})$. As such, all samples from $q^*(\boldsymbol{\theta})$ will lie within a mode of $f(\boldsymbol{\theta}|\mathbf{x})$. Belief propagation is motivated by reversing the KL divergence that drives mean field variational Bayes, i.e.

$$q^*(\boldsymbol{\theta}) = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \left(f(\boldsymbol{\theta}|\mathbf{x}) \parallel \prod_{i=1}^m q(\boldsymbol{\theta}_i) \right). \quad (1.22)$$

By reversing the direction of KL divergence, $q^*(\boldsymbol{\theta})$ is selected by maximising the probability of $f(\boldsymbol{\theta}|\mathbf{x})$ being under $q(\boldsymbol{\theta}_i)$. As such, the optimal approximate distribution $q^*(\boldsymbol{\theta})$ will cover all modes of $f(\boldsymbol{\theta}|\mathbf{x})$. The difference in the approximation resulting from the different directions of KL divergence is illustrated in Figure 1.7. The mean seeking approximation that arises can be solved by a moment matching problem. However, this problem has the potential to become unboundedly complex. As mentioned before, one way to control the complexity of projections required is to restrict approximating factors to be in an exponential family. Both ADF and EP use the same direction of KL divergence as belief propagation, but constrain the approximating distributions to be from the exponential family to aid their implementation.

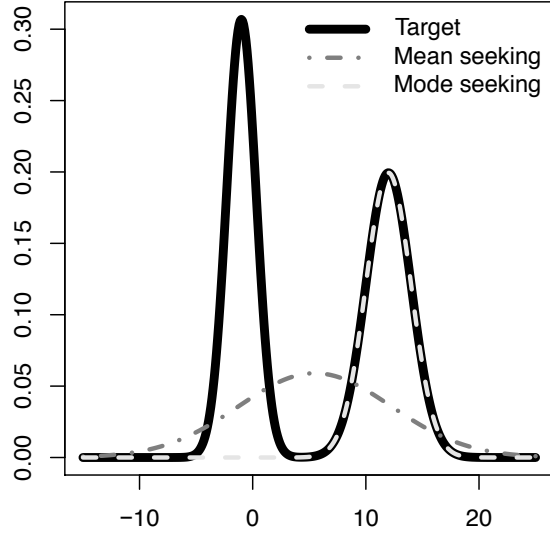


Figure 1.6: *Plot comparing density estimates using both directions of the KL divergence. The approximating density was Gaussian and the target density was a mixture of two Gaussians.*

1.12.3 Expectation propagation

We now explain the general EP schematic, re-iterating points previously made in this section. For m independent data points \mathbf{x}_i where $i \in \{1, \dots, m\}$ and latent variable \mathbf{u} consider an intractable posterior

$$\begin{aligned} p(\mathbf{u}|\mathbf{x}) &= \frac{p_0(\mathbf{u}) \prod_{i=1}^m p(\mathbf{x}_i|\mathbf{u})}{\int p_0(\mathbf{u}) \prod_{i=1}^m p(\mathbf{x}_i|\mathbf{u}) d\mathbf{u}} \\ &= Z^{-1} p_0(\mathbf{u}) \prod_{i=1}^m p(\mathbf{x}_i|\mathbf{u}), \end{aligned} \quad (1.23)$$

where

$$p(\mathbf{x}_i|\mathbf{u}) = f_i(\mathbf{u}), \quad p_0(\mathbf{u}) = f_0(\mathbf{u}), \quad Z = \int p_0(\mathbf{u}) \prod_{i=1}^m p(\mathbf{x}_i|\mathbf{u}) d\mathbf{u},$$

and $f_0(\mathbf{u})$ is a member of the exponential family. Note the posterior in equation (1.23) is proportional to

$$p(\mathbf{u}|\mathbf{x}) \propto f_0(\mathbf{u}) \prod_{i=1}^m f_i(\mathbf{u}).$$

We wish to obtain a global approximation of the posterior with a density q selected

from the exponential family (denoted by Q)

$$q(\mathbf{u}) \propto \prod_{i=0}^m q_i(\mathbf{u}),$$

where $q_0(\mathbf{u}) = f_0(\mathbf{u})$. The optimal global approximation is the density which minimises the KL divergence to the posterior

$$q^*(\mathbf{u}) = \operatorname{argmin}_{q \in Q} \operatorname{KL}(p(\mathbf{u}|\mathbf{x})||q(\mathbf{u})).$$

Since the approximating density family is selected to be exponential, finding the optimal approximation reduces to a simple moment matching problem. This moment matching problem is not feasible over the whole posterior, so we instead conduct it on each site. The product of these site approximations is equivalent to the approximation over the whole posterior. As such, our goal is to find an approximation $q_i(\mathbf{u})$ which minimises the KL divergence over a site of the likelihood factor $p(\mathbf{x}_i|\mathbf{u})$. In other words, we require the KL projection of a site of our target distribution onto the approximating density family. However, as the target distribution is outside the exponential family space, we instead use a tractable hybrid distribution $h_i(\mathbf{u})$ that lies in between the target distribution and the approximating distribution, as shown in Figure 1.7. Each hybrid distribution (also

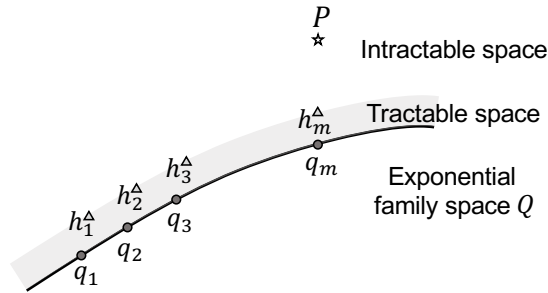


Figure 1.7: *Visualisation of the EP process (Barthelme, 2016^[2]). The hybrid approximations are marked with triangles, the approximations are marked with circles and the target density is marked with a star.*

called tilted distribution) is found by removing an approximating site from the global approximation (removing the i th site forms a cavity distribution denoted by $q_{-i}(\mathbf{u})$) and replacing it with the equivalent true likelihood site. Formally the i th hybrid is

$$h_i(\mathbf{u}) = f_i(\mathbf{u})q_{-i}(\mathbf{u}),$$

where the i th cavity distribution $q_{-i}(\mathbf{u})$ is given by

$$q_{-i}(\mathbf{u}) = \prod_{i' \neq i}^m q_{i'}(\mathbf{u}).$$

We can now project this hybrid to the approximating family by computing the moments, which results in a new approximation of the whole posterior with features of the hybrid

$$q^{\text{new}}(\mathbf{u}) = \text{proj}[h_i(\mathbf{u})].$$

In the special case where the approximating family is Gaussian this problem simplifies to finding the mean, variance and normalising factor of the hybrid distribution. Using the new global approximation of the posterior we must now update the site $q_i(\mathbf{u})$ such that $q_i(\mathbf{u})q_{-i}(\mathbf{u})$ has the same moments as the hybrid distribution, i.e.

$$q_i^{\text{new}}(\mathbf{u}) = \frac{q^{\text{new}}(\mathbf{u})}{q_{-i}(\mathbf{u})}.$$

We hope that by iterating projections of the hybrid distributions onto the exponential family space and updating each $q_i(\mathbf{u})$, we find a stationary point which minimises the KL distance to the target density, as shown in Figure [1.7](#). At this point a global approximation to the posterior can be obtained by taking the product of the partitions,

$$q(\mathbf{u}) = \prod_{i=0}^m q_i(\mathbf{u}).$$

This process forms the EP algorithm, and is presented explicitly in Algorithm [1](#). By using natural parameters further simplification is available, where the optimal parameters are found by linear algebra.

Algorithm 1 *A generalised version of the EP algorithm.*

Initialise: by setting $q(\mathbf{u}) = \prod_{i=0}^m q_i(\mathbf{u})$, $q_0(\mathbf{u}) = f_0(\mathbf{u})$ and for $1 \leq i \leq m$, $q_i(\mathbf{u}) = 1$.

Cycle: Pick $i = 1, \dots, m$:

Remove $q_i(\mathbf{u})$ from $q(\mathbf{u})$ to get the cavity distribution $q_{-i}(\mathbf{u})$,

$$q_{-i}(\mathbf{u}) \propto \prod_{i' \neq i}^m q_{i'}(\mathbf{u}). \quad (1.24)$$

Update the tilted distribution by replacing the approximating likelihood $q_i(\mathbf{u})$ by the exact one $f_i(\mathbf{u})$,

$$h_i(\mathbf{u}) = f_i(\mathbf{u})q_{-i}(\mathbf{u}).$$

Project $h_i(\mathbf{u})$ back to the exponential family

$$q^{\text{new}}(\mathbf{u}) = \text{proj}[h_i(\mathbf{u})].$$

Update the approximate terms

$$q_i^{\text{new}}(\mathbf{u}) \propto \frac{q^{\text{new}}(\mathbf{u})}{q_{-i}(\mathbf{u})} \quad (1.25)$$

until all q_i converge.

After convergence is reached obtain a tractable approximation from

$$q(\mathbf{u}) = \prod_{i=0}^m q_i(\mathbf{u}).$$

1.12.4 Message passing

Minka (2005)^[44] streamlines EP and other variational inference algorithms into a framework known as message passing, which allows significant algebraic and computational simplifications for large factor graphs. In this section we focus on the message passing approach for simple EP models and refer interested readers to Minka (2005)^[44] for other variational methods.

A message is simply a function defined by a factor node that takes a subset of parental stochastic nodes as an input. The EP problem of obtaining a KL projection onto the required exponential family can be solved in terms of messages passed between neighbouring nodes of a relevant factor graph. Since the approximating family is constrained to be exponential, the messages can be expressed in natural parameters.

Let f_i , $1 \leq i \leq m$ denote factor nodes and \mathbf{u} denote a stochastic node that corresponds to a latent variable. For the models presented in this thesis, the neighbours of any i th factor node are the stochastic node \mathbf{u} . Kim & Wand (2016)^[31] provide additional details for larger models. Given the latent variable \mathbf{u} and observed data \mathbf{x} the joint density function can be written as

$$p(\mathbf{u}, \mathbf{x}) = \prod_{i=1}^m f_i(\mathbf{u}).$$

For the simple factor graphs presented in this thesis, the message passing EP algorithm can be reduced to updating messages from the stochastic node to factor node as

$$m_{\mathbf{u} \rightarrow f_i}(\mathbf{u}) \leftarrow \prod_{i' \neq i}^m m_{f_{i'} \rightarrow \mathbf{u}}(\mathbf{u}), \quad (1.26)$$

and updating the message from the factor to stochastic node as

$$m_{f_i \rightarrow \mathbf{u}}(\mathbf{u}) \leftarrow \frac{\text{proj}[Z^{-1} m_{\mathbf{u} \rightarrow f_i}(\mathbf{u}) f_i(\mathbf{u})]}{m_{\mathbf{u} \rightarrow f_i}(\mathbf{u})}, \quad (1.27)$$

where Z is the normalising factor. In terms of the Algorithm [1](#), equation [\(1.26\)](#) for message passing EP corresponds to forming the cavity distribution in equation [\(1.24\)](#), and equation [\(1.27\)](#) corresponds updating the $q_i(\mathbf{u})$ approximation as in equation [\(1.25\)](#).

Once the messages have reached convergence, the KL optimal q -densities are obtained via

$$q(\mathbf{u}) \propto \prod_{i=1}^m m_{f_i \rightarrow \mathbf{u}}(\mathbf{u}).$$

This process is presented consisely in Algorithm [2](#). Kim & Wand (2016)^[31] suggest using relative change in the approximate marginal log-likelihood as a stopping criterion. We refer readers there for further details.

Algorithm 2 *A generalised version of the message passing EP algorithm.*

Initialise: by setting $m_{f_i \rightarrow \mathbf{u}}(\mathbf{u}) = 1$, where $1 \leq i \leq m$.

Cycle: Pick $i = 1, \dots, m$:

 Get the cavity distribution,

$$m_{\mathbf{u} \rightarrow f_i}(\mathbf{u}) \leftarrow \prod_{i' \neq i}^m m_{f_{i'} \rightarrow \mathbf{u}}(\mathbf{u}).$$

 Project the tilted distribution onto the exponential family and update the approximate terms

$$m_{f_i \rightarrow \mathbf{u}}(\mathbf{u}) \leftarrow \frac{\text{proj}[Z^{-1} m_{\mathbf{u} \rightarrow f_i}(\mathbf{u}) f_i(\mathbf{u})]}{m_{\mathbf{u} \rightarrow f_i}(\mathbf{u})}$$

 until all $m_{f_i \rightarrow \mathbf{u}}$ converge.

After convergence is reached obtain a tractable approximation from

$$q(\mathbf{u}) = \prod_{i=1}^m m_{f_i \rightarrow \mathbf{u}}(\mathbf{u}).$$

1.13 Thesis structure

We have now explained the background information required and from this point on present novel research. Chapter 2 explores a simple model for the probit link where we estimate only the variance parameter σ^2 i.e. (the model only has a random intercept). Chapter 3 extends the work in Chapter 2 to the general case where there are multiple fixed and random effects. Chapter 4 applies the work of Chapters 2 and 3 to the logistic link function. Chapter 5 continues from Chapter 3 by applying the work from the previous chapters to count data with the Poisson link and the negative binomial link. Chapter 6 explores models with crossed random effects. Chapter 7 applies the methodology implemented in this thesis to two real datasets, one provided by the Australian Red Cross Blood Service. The thesis ends with a discussion and concluding remarks in Chapter 8.

Chapter 2

Expectation propagation for the simplest one level probit mixed model

In this chapter we develop methodology for frequentist inference of binary probit GLMMs for the simplest random intercepts only model. Our goal is to approximate the maximum likelihood of the parameter for variance between groups only (i.e. random intercepts only model) with 95% confidence intervals. Since this model assumes only one source of variance, a change in notation occurs from the covariance matrix Σ shown in Section 1.7 to the scalar variance σ^2 . This model provides a good starting point for testing new methodology since only one dimensional integrals are required, as opposed to the integrals required for more complex models. Additionally, any difficulties caused by matrices are negated. We trial various methods of calculating confidence intervals and show how to obtain best predictors.

In this chapter, we assume a balanced dataset, where all m groups have the same n number of observations in them. For observed values of

$$y_{ij}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

where $y_{ij} \in \{0, 1\}$, the probit binary mixed model form is

$$y_{ij}|u_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\Phi(u_i)), \quad u_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), \quad (2.1)$$

where u_i is a scalar unobserved latent variable. The variance parameter likelihood can

be expressed as

$$\ell(\sigma^2) = \sum_{i=1}^m \ell_i(\sigma^2),$$

where

$$\ell_i(\sigma^2) \equiv \log \int_{-\infty}^{\infty} \prod_{j=1}^n \Phi((2y_{ij} - 1)u_i) (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/2\sigma^2) du_i, \quad (2.2)$$

and the best predictor of u_i is

$$\text{BP}(u_i) = \frac{\int_{-\infty}^{\infty} u_i \prod_{j=1}^n \Phi((2y_{ij} - 1)u_i) \exp(-u_i^2/2\sigma^2) du_i}{\int_{-\infty}^{\infty} \prod_{j=1}^n \Phi((2y_{ij} - 1)u_i) \exp(-u_i^2/2\sigma^2) du_i}.$$

We denote the maximum likelihood estimate of σ^2 by

$$\widehat{\sigma^2} = \underset{\sigma^2}{\operatorname{argmax}} \ell(\sigma^2).$$

Calculation of the maximum likelihood estimator and best predictor are complicated by the intractable integral arising in equation (2.2). Each $\ell_i(\sigma^2)$ can be approximated and summed to obtain the full log-likelihood. Traditionally these univariate integrals are solved using quadrature. We develop an EP scheme for estimating the likelihood surface by approximating each $\ell_i(\sigma^2)$ and summing them to obtain the full log-likelihood. We compare its performance to quadrature. Our aim is to show both methods provide reasonable and similar estimates.

Details of the quadrature approach to estimating the likelihood surface are given in Section 2.1. Our novel method using EP is explained in Section 2.2. Section 2.4 explores different ways to obtain the maximum likelihood estimate and confidence intervals for both our novel method and the quadrature approach. Details on obtaining best predictors for this model are given in Section 2.5, before comparing the approaches of likelihood estimation in a simulation study in Section 2.6.

2.1 Traditional quadrature likelihood approximation

Likelihood surface approximation by quadrature is easy to implement using software such as the R function `integrate()` in the “stats” package (R Core Team, 2019^[56]) as discussed in Section 1.11.2. However, one must be careful to avoid issues with numerical instability. As such, we express the integrand in equation (2.2) as a function that attains

a maximum value of 1 over $u \in \mathbb{R}$, where the i th integral arising in each $\ell_i(\sigma^2)$ is

$$I_i(\sigma^2) = \exp \left\{ \max_{u_0 \in \mathbb{R}} (h_i(u_0)) \right\} \int_{-\infty}^{\infty} \exp \left\{ h_i(u) - \max_{u_0 \in \mathbb{R}} (h_i(u_0)) \right\} du$$

and

$$h_i(u) \equiv \sum_{j=1}^n \log \Phi((2y_{ij} - 1)u) - \frac{u^2}{2\sigma^2}.$$

The first and second derivatives of $h_i(u)$ are

$$h_i'(u) = \sum_{j=1}^n \zeta'((2y_{ij} - 1)u)(2y_{ij} - 1) - \frac{u}{\sigma^2}$$

and

$$h_i''(u) = \sum_{j=1}^n \zeta''((2y_{ij} - 1)u) - \frac{1}{\sigma^2}.$$

Assuming $\zeta''(x) < 0$ for all $x \in \mathbb{R}$, then $h_i''(u) < 0$ for all $u \in \mathbb{R}$. From this fact,

- h_i is a strictly concave function over \mathbb{R}
- $\exp(h_i(u))$ is log-concave
- $\lim_{u \rightarrow -\infty} h_i'(u) = -\infty$ and $\lim_{u \rightarrow \infty} h_i'(u) = +\infty$.

As such, each $\ell_i(\sigma^2)$ of the likelihood is calculated as

$$\frac{1}{2} \log(2\pi\sigma^2) + \ell_i(\sigma^2) = h(u_{0i}) + \log \int_{-\infty}^{\infty} \exp(h_i(u_i) - h_i(u_{0i})) du_i,$$

where $h_i'(u_{0i}) = 0$ and the unique root u_{0i} is found using a bisection search, where the starting values are selected -1 and 1 to be for the lower and upper bounds respectively.

2.2 Expectation propagation likelihood approximation

We now propose a novel approach to likelihood approximation using EP (the log-likelihood is denoted by $\ell(\sigma^2)$) which is more amenable to cases involving higher dimensional integrals than traditional quadrature routines. As discussed in Section [1.12.3](#), the EP approximation is motivated by minimisation of a KL divergence criterion (see equation [\(1.19\)](#)), which selects an unnormalised normal density function to replace

each

$$\Phi((2y_{ij} - 1)u_i), \quad 1 \leq j \leq n$$

in equation (2.2). Subsequently, the integrand is proportional to a product of univariate normal density functions which have explicit forms in probit case. Furthermore, the calculation of $\ell(\sigma^2)$ requires only fixed point iteration, so there is no need for any numerical integration. A major downfall of EP implementation is the high algebraic overhead. We aim to minimise this by using message passing similar to Kim & Wand (2017).³²

Consider the family of unnormalised normal density functions written in exponential family form,

$$f_{\text{UN}}(x) = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \eta_0 \\ \eta_1 \\ \eta_2 \end{bmatrix} \right\} \quad (2.3)$$

with natural parameters $\eta_0, \eta_1 \in \mathbb{R}$ and $\eta_2 < 0$. The goal of the EP problem is to find the optimal natural parameters, η_0^*, η_1^* and η_2^* , which minimise $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$ where $f_{\text{input}} \in L_1$. This solution is referred to as the KL projection onto the family of unnormalised normal density functions, and is written as

$$\text{proj}_{\text{UN}}[f_{\text{input}}](x) = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta}^* \right\}, \quad (2.4)$$

where

$$\boldsymbol{\eta}^* \equiv \begin{bmatrix} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{bmatrix},$$

$$(\eta_0^*, \eta_1^*, \eta_2^*) = \underset{(\eta_0, \eta_1, \eta_2) \in H}{\text{argmin}} \text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$$

and H is the set of all allowable natural parameters. In the special case of KL projection onto the unnormalised univariate normal family this problem simplifies further to moment-matching, where $(\eta_0^*, \eta_1^*, \eta_2^*)$ is the unique vector that matches the zeroth, first and second order moments of f_{UN} and f_{input} . For the case of probit binary GLMMs,

EP requires repeated KL projections of the form

$$f_{\text{input}}(x) = \Phi(c_0 + c_1 x) \exp(\eta_1^{\text{input}} x + \eta_2^{\text{input}} x^2) \quad (2.5)$$

onto an unnormalised normal distribution, where $c_0 = 0$, $c_1 = 2y_{ij} - 1$, $x = u_i$, $\eta_1^{\text{input}} \in \mathbb{R}$ and $\eta_2^{\text{input}} < 0$. As such, we seek $\boldsymbol{\eta}^*$ such that

$$\int_{-\infty}^{\infty} x^k \Phi(c_0 + c_1 x) \exp\left\{\left[\begin{array}{c} x \\ x^2 \end{array}\right]^{\top} \boldsymbol{\eta}^{\text{input}}\right\} dx = \int_{-\infty}^{\infty} x^k \exp\left\{\left[\begin{array}{c} 1 \\ x \\ x^2 \end{array}\right]^{\top} \boldsymbol{\eta}^*\right\} dx. \quad (2.6)$$

Consider Result [3](#)

Result 3. For an unnormalised input function $f \in L_1$ such that $f \geq 0$ for all $x \in \mathbb{R}$, where $C_f \equiv \int_{\mathbb{R}} f(x) dx$, the projection onto the unnormalised normal family is

$$\text{proj}_{UN}[f] = C_f \text{proj}_N[f/C_f](x),$$

where proj_N is the projection onto the normal family.

By Result [3](#), obtaining the natural parameters $\boldsymbol{\eta}^*$ for projection onto the unnormalised normal family follows from obtaining the projection onto the normal family. More explicitly, the optimal natural parameters, η_1^* and η_2^* , are given according to the projection of the normalised function $f_{\text{input}}/C_{f_{\text{input}}}$ onto the normal family. We can subsequently use these optimal natural parameters to find the normalising natural parameter η_0^* via Result [4](#) and thus obtain the projection onto unnormalised normal family.

Result 4. When the input density follows the form of equation [\(2.5\)](#), η_0^* is given by

$$\begin{aligned} \eta_0^* &= \log(C_{f_{\text{input}}}) - A(\eta_1^*, \eta_2^*) - \frac{1}{2} \log(2\pi), \\ &= \log(C_{f_{\text{input}}}) + (\eta_1^*)^2 / (4\eta_2^*) + \frac{1}{2} \log(-\eta_2^*/\pi). \end{aligned}$$

Thus to obtain the required projection, we first obtain the optimal natural parameters η_1^* and η_2^* to project onto the normal family as is presented in Result [5](#).

Result 5. Given f follows the form of equation (2.5), the projection onto the univariate normal family is given by

$$\text{proj}_N[f] = \exp(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*))h(x)$$

where

$$\boldsymbol{\eta}_{-1}^{\text{input}} \equiv \begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix} = k_{\text{probit}}(\boldsymbol{\eta}_{-1}^{\text{input}}; c_0, c_1),$$

$k_{\text{probit}}\left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1\right)$ is defined in Definition 10 and $\mathbf{T}(x)$ and $h(x)$ follow from Section 1.5.2.1

Definition 10. For primary scalar arguments $a_1 \in \mathbb{R}$ and $a_2 < 0$ and auxiliary scalar arguments $c_0, c_1 \in \mathbb{R}$, the function $k_{\text{probit}} : H \rightarrow H$ is given by

$$k_{\text{probit}}\left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1\right) = \begin{bmatrix} r_5(a_1 + r_3c_1) \\ r_5a_2 \end{bmatrix},$$

with

$$r_1 = \sqrt{2(2 - c_1^2 a_2^{-1})}, \quad r_2 = (2c_0 - c_1 a_2^{-1} a_1) r_1^{-1}, \quad r_3 = 2\zeta'(r_2) r_1^{-1},$$

$$r_4 = -2\zeta''(r_2) r_1^{-2} \quad \text{and} \quad r_5 = (a_2 + r_4 c_1^2)^{-1} a_2.$$

Using Result 5 we now obtain the normalising natural parameter η_0^* to find the projection onto unnormalised normal family.

2.2.1 Projection onto the unnormalised normal family

Recalling the moment matching problem described by equation (2.6) and Result 4, the normalising factor can be shown to be

$$C_{f_{\text{input}}} = \int_{\mathbb{R}} f_{\text{input}}(x) dx = (2\pi)^{-1/2} \exp(A(\boldsymbol{\eta}^{\text{input}})) \Phi(r_2),$$

where r_2 is given in Definition 10 and $A(\boldsymbol{\eta})$ is defined in Section 1.5. By Result 4

$$\eta_0^* = \log \Phi(r_2^{\text{input}}) + \frac{1}{4}(\eta_1^*)^2/\eta_2 - \frac{1}{4}(\eta_1^{\text{input}})^2/\eta_2^{\text{input}} + \frac{1}{2}\log(\eta_2^*/\eta_2^{\text{input}}).$$

To obtain η_0^* we introduce the $c_{\text{probit}}(\mathbf{a}, \mathbf{b}; c_0, c_1)$ function in Definition 11, for $\mathbf{a} = [a_1 \ a_2]^\top$ and $\mathbf{b} = [b_1 \ b_2]^\top$.

Definition 11. Consider first, primary scalar arguments a_1, a_2, b_1 and b_2 , and auxiliary scalar arguments c_0 and c_1 . Then the function $c_{\text{probit}} : H \rightarrow \mathbb{R}$ is given by

$$c_{\text{probit}} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}; c_0, c_1 \right) \equiv \log \Phi(r_2) + \frac{1}{4}b_1^2/b_2 - \frac{1}{4}a_1^2/a_2 + \frac{1}{2} \log(b_2/a_2),$$

where r_2 follows from Definition 10.

To summarise, the projection of the input function onto the unnormalised normal family is obtained as in Result 6.

Result 6. For an unnormalised input function of the form of equation (2.5),

$$\text{proj}_{UN}[f_{\text{input}}] = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{bmatrix} \right\},$$

where

$$\eta_0^* = c_{\text{probit}} \left(\begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}, \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix}; c_0, c_1 \right)$$

and

$$\begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix} = k_{\text{probit}} \left(\begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}; c_0, c_1 \right).$$

Note the forms of these functions are useful in the next section where we show how to organise the projections required.

2.2.2 Message passing formulation

We now express the EP approximation of $\ell(\sigma^2)$ using message passing updates as presented in Kim & Wand (2017).^[68] Note each $\ell_i(\sigma^2)$ can be written as

$$\ell_i(\sigma^2) = \log \int_{-\infty}^{\infty} \left(\prod_{j=1}^n p(y_{ij}|u_i) \right) p(u_i; \sigma^2) du_i, \quad (2.7)$$

where

$$p(y_{ij}|u_i) \equiv \Phi((2y_{ij} - 1)u_i) \quad \text{and} \quad p(u_i; \sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/(2\sigma^2))$$

are respectively the conditional density function of each response given its random effect and the density function of the random effect. The alternate expression

$$p(u_i; \sigma^2) = \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \eta_{\sigma^2} \right\}, \quad \text{where} \quad \eta_{\sigma^2} \equiv \begin{bmatrix} -\frac{1}{2} \log(2\pi\sigma^2) \\ 0 \\ -1/(2\sigma^2) \end{bmatrix} \quad (2.8)$$

is more amenable to the message passing approach to follow and is also worth noting. Using factor graphs (discussed in Section 1.6.2), Figure 2.1 visualises the dependence structure of the product in equation (2.7), where the circular stochastic node corresponds to the random vector u_i and solid squares indicate the $n+1$ factor nodes. The dependence of each factor node on the stochastic node u_i is demonstrated through the edges.

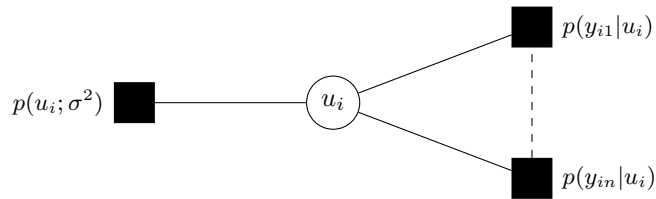


Figure 2.1: Factor graph representation of the product structure of the integrand in equation (2.7).

With this factor graph in place implementing EP is analogous to the Bayesian approach in Minka (2005).^[44] As previously explained, the EP approximation of $\ell_i(\sigma^2)$ is motivated by minimization of a KL divergence and imposition of exponential family

constraints. Consider that

$$\tilde{p}(y_{ij}|u_i) = \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \boldsymbol{\eta}_{ij} \right\}, \quad 1 \leq j \leq n$$

is initialised such that u_i is an unnormalised normal density function. Then for each $j = 1, \dots, n$, the $\boldsymbol{\eta}_{ij}$ update requires minimisation of

$$\text{KL} \left(p(y_{ij}|u_i) \left(\prod_{j' \neq j} \tilde{p}(y_{ij'}|u_i) \right) p(u_i; \sigma^2) \parallel \left(\prod_{j'=1}^n \tilde{p}(y_{ij'}|u_i) \right) p(u_i; \sigma^2) \right) \quad (2.9)$$

as functions of u_i . Thus we can use Result [6](#) to perform updates until convergence of $\boldsymbol{\eta}_{ij}$ s. Using Figure [2.1](#), EP can be compartmentalised by the notion of message passing as shown in Section 4.1 of Minka (2005).^{[44](#)} First let us define messages passed from the factor $p(y_{ij}|u_i)$ to u_i as

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \equiv \tilde{p}(y_{ij}|u_i).$$

Messages from the factor $p(y_{ij}|u_i)$ to the stochastic node u_i are updated by equation (60) of Minka (2005),^{[44](#)} where $\alpha = 1$ and $s' = 1$ (since the KL divergence we are working with is unnormalised). Simplifications given we have one stochastic node result in the following expression for equation [\(2.9\)](#),

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}_{\text{UN}}[m_{u_i \rightarrow p(y_{ij}|u_i)}(u_i) p(y_{ij}|u_i)](u_i)}{m_{u_i \rightarrow p(y_{ij}|u_i)}(u_i)}, \quad 1 \leq j \leq n. \quad (2.10)$$

Similarly the update of message passed from $p(u_i; \sigma^2)$ to u_i is

$$m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}_{\text{UN}}[m_{u_i \rightarrow p(u_i; \sigma^2)}(u_i) p(u_i; \sigma^2)](u_i)}{m_{u_i \rightarrow p(u_i; \sigma^2)}(u_i)}. \quad (2.11)$$

By equation (54) of Minka (2005)^{[44](#)} the updates of stochastic node to factor messages are

$$m_{u_i \rightarrow p(y_{ij}|u_i)}(u_i) = m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) \prod_{j' \neq j} m_{p(y_{ij'}|u_i) \rightarrow u_i}(u_i), \quad 1 \leq j \leq n \quad (2.12)$$

and

$$m_{u_i \rightarrow p(u_i; \sigma^2)}(u_i) = \prod_{j=1}^n m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i). \quad (2.13)$$

We now seek any algebraic simplifications of the key messages. Recall that $p(u_i; \sigma^2)$ can be written using natural parameters in the form of equation (2.8) and that the unnormalised normal density constraint is enforced on equations (2.10) and (2.12). Then

$$m_{u_i \rightarrow p(u_i; \sigma^2)}(u_i) = \exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{u_i \rightarrow p(u_i; \sigma^2)} \end{array} \right\}. \quad (2.14)$$

Substituting the above forms into equation (2.11) leads to

$$\begin{aligned} m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) &\leftarrow \frac{\text{proj}_{\text{UN}} \left[\exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{u_i \rightarrow p(u_i; \sigma^2)} \end{array} \right\} \exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{\sigma^2} \end{array} \right\} \right]}{\exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{u_i \rightarrow p(u_i; \sigma^2)} \end{array} \right\}} \\ &= \frac{\exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{u_i \rightarrow p(u_i; \sigma^2)} \end{array} \right\} \exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{\sigma^2} \end{array} \right\}}{\exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{u_i \rightarrow p(u_i; \sigma^2)} \end{array} \right\}} \\ &= \exp \left\{ \begin{array}{c} \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \\ \boldsymbol{\eta}_{\sigma^2} \end{array} \right\}. \end{aligned}$$

This implies that the message $m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) = p(u_i; \sigma^2)$ is constant throughout the message passing updates. As such, we now set

$$\boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} \leftarrow \boldsymbol{\eta}_{\sigma^2}. \quad (2.15)$$

For convenience,

$$\boldsymbol{\eta}^\otimes \equiv \boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)}.$$

Recall that from equation (2.14)

$$m_{u_i \rightarrow p(y_{ij}|u_i)}(u_i) = \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \boldsymbol{\eta}^\otimes \right\} = \exp(\eta_0^\otimes) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2).$$

Substituting this into equation (2.10) leads to

$$\begin{aligned} m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) &\leftarrow \frac{\text{proj}_{\text{UN}} \left[\exp(\eta_0^\otimes) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2) \Phi((2y_{ij} - 1)u_i) \right]}{\exp(\eta_0^\otimes) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)} \\ &= \frac{\exp(\eta_0^\otimes) \text{proj}_{\text{UN}} \left[\Phi((2y_{ij} - 1)u_i) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2) \right]}{\exp(\eta_0^\otimes) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)} \\ &= \frac{\text{proj}_{\text{UN}} \left[\Phi(c_0 + c_{1_{ij}} u_i) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2) \right]}{\exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)}, \end{aligned}$$

where $c_0 = 0$ and $c_{1_{ij}} \equiv 2y_{ij} - 1$. By utilising Result 6,

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \right\}, \quad (2.16)$$

where the linear and quadratic coefficient updates are

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} \leftarrow k_{\text{probit}}(\boldsymbol{\eta}_{1:2}^\otimes; c_0, c_{1_{ij}}) - \boldsymbol{\eta}_{1:2}^\otimes \quad (2.17)$$

and the constant coefficient update is

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 \leftarrow c_{\text{probit}}(\boldsymbol{\eta}_{1:2}^\otimes; (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_0, c_{1_{ij}}). \quad (2.18)$$

Using the simplification of equation (2.10) and (2.11), equation (2.12) can be shown to

be

$$m_{u_i \rightarrow p(y_{ij}|u_i)}(u_i) \leftarrow \exp \left\{ \left[\begin{array}{c} 1 \\ u_i \\ u_i^2 \end{array} \right]^\top \boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)} \right\}, \quad (2.19)$$

where

$$\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)} \leftarrow \boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} + \sum_{j' \neq j} \boldsymbol{\eta}_{p(y_{ij'}|u_i) \rightarrow u_i}.$$

The scheme of the message passing approach to EP can be summarised as in Section 6 of Minka (2005):⁴⁴

1. Initialise all factor to stochastic node messages.
2. Cycle until convergence of all factor to stochastic node message:

For each factor:

- (a) Compute the messages passed to the factor via equation (2.14) or equation (2.19).
- (b) Compute the messages passed from the factor via equation (2.15) or equation (2.16).

The EP approximation to each log-likelihood $\ell_i(\sigma^2)$ component is given by

$$\ell_i(\sigma^2) = \log \int_{\mathbb{R}} \left(\prod_{j=1}^n m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \right) m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) du_i. \quad (2.20)$$

The success of EP depends on each of the messages in equation (2.20) being an unnormalised normal density. This allows results in a closed form solution to the

integral as follows:

$$\begin{aligned}
& \int_{\mathbb{R}} \left(\prod_{j=1}^n m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \right) m_{p(u_i; \sigma^2) \rightarrow u_i}(u_i) du_i \\
&= \int_{\mathbb{R}} \prod_{j=1}^n \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \right\} \exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} \right\} du_i \\
&= (2\pi)^{1/2} \exp \left\{ \left(\boldsymbol{\eta}_{\sigma^2} + \sum_{j=1}^n \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \right)_0 + A \left(\left(\boldsymbol{\eta}_{\sigma^2} + \sum_{j=1}^n \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \right)_{1:2} \right) \right\}.
\end{aligned}$$

The full algorithm for the approximation of $\ell(\sigma^2)$ using EP is provided in Algorithm [3](#).

Algorithm 3 *Explicit expression of the algorithm used for the message passing approach to EP in the random intercepts only model.*

Inputs: y_{ij} , $\mathbf{x}_{ij}^{\mathbf{R}}$, $1 \leq i \leq m$, $1 \leq j \leq n$; $\Sigma(d^{\mathbf{R}} \times d^{\mathbf{R}}$, is symmetric and positive definite).
 Set constants: $c_0 \leftarrow 0$, $c_{1_{ij}} \leftarrow 2y_{ij} - 1$; $1 \leq i \leq m$, $1 \leq j \leq n$,

$$\boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} \leftarrow \boldsymbol{\eta}_{\sigma^2} \equiv \begin{bmatrix} -\frac{1}{2} \log(2\pi\sigma^2) \\ 0 \\ 1/(2\sigma^2) \end{bmatrix}, \quad 1 \leq i \leq m. \quad (2.21)$$

For $i = 1, \dots, m$:

Initialise: $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \in \mathbb{R}$, $1 \leq j \leq n$ as per equation (2.25).

Cycle:

$$\text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) \leftarrow \sum_{j=1}^n \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$$

For $j = 1, \dots, n$:

$$\begin{aligned} \boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)} &\leftarrow \boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) - \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \\ (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} &\leftarrow k_{\text{probit}} \left((\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}; c_0, c_{1_{ij}} \right) \\ &\quad - (\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2} \end{aligned}$$

until all natural parameter vectors converge.

For $j = 1, \dots, n$:

$$\begin{aligned} (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 &\leftarrow c_{\text{probit}} \left((\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}, \right. \\ &\quad \left. (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} + (\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}; c_0, c_{1_{ij}} \right). \end{aligned}$$

$$\text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) \leftarrow \sum_{j=1}^n \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$$

Output: The full approximate log-likelihood is given by

$$\begin{aligned} \underline{\ell}(\sigma^2) &= (m/2) \log(2\pi) + \sum_{i=1}^m \left\{ \left(\boldsymbol{\eta}_{\sigma^2} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) \right)_0 \right. \\ &\quad \left. + A \left\{ \left(\boldsymbol{\eta}_{\sigma^2} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) \right)_{1:2} \right\} \right\}, \end{aligned}$$

where $A(\boldsymbol{\eta})$ is defined in equation (1.5) and $\boldsymbol{\eta}_{\sigma^2}$ follows from equation (2.21).

2.2.3 Starting values for Algorithm 3

The EP message passing algorithm proposed relies on good starting values for convergence. We now derive starting values for $\eta_{p(y_{ij}|u_i) \rightarrow u_i}$ using a Taylor series expansion. Note that

$$\log p(y_{ij}|u_i) = \zeta(a_{ij}) - \log(2),$$

where $a_{ij} \equiv (2y_{ij} - 1)u_i$ and ζ is defined in Section 1.5.2.1. Let \hat{u}_i be a Laplace approximation to u_i . Now consider the following Taylor series expansion of the data dependent component of $\ell(\sigma^2)$:

$$\begin{aligned} \zeta(a_{ij}) &= \zeta(\hat{a}_{ij} + (u_i - \hat{u}_i)(2y_{ij} - 1)) \\ &= \zeta(\hat{a}_{ij}) + (u_i - \hat{u}_i)(2y_{ij} - 1)\zeta'(\hat{a}_{ij}) + \frac{1}{2}((u_i - \hat{u}_i)(2y_{ij} - 1))^2\zeta''(\hat{a}_{ij}) + \dots \\ &= \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\eta}_{ij} + \dots, \end{aligned}$$

where $\hat{a}_{ij} \equiv (2y_{ij} - 1)\hat{u}_i$ and

$$\check{\eta}_{ij} = \begin{bmatrix} \zeta(\hat{a}_{ij}) \\ (2y_{ij} - 1)\zeta'(\hat{a}_{ij}) \\ \frac{1}{2}\zeta''(\hat{a}_{ij}) \end{bmatrix}.$$

It follows that the quadratic approximation to $\log p(y_{ij}|u_i)$ based on Taylor expansion about \hat{u}_i is $\log \check{p}(y_{ij}|u_i)$ where

$$\check{p}(y_{ij}|u_i) \equiv \exp \left\{ \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\eta}_{ij} \right\}. \quad (2.22)$$

The starting value recommendation for $\eta_{p(y_{ij}|u_i)}$ is based on replacement of $p(y_{ij}|u_i)$ by $\check{p}(y_{ij}|u_i)$ in equation (2.15):

$$m_{\check{p}(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}[m_{u_i \rightarrow \check{p}(y_{ij}|u_i)}(u_i)\check{p}(y_{ij}|u_i)](u_i)}{m_{u_i \rightarrow \check{p}(y_{ij}|u_i)}(u_i)} = \check{p}(y_{ij}|u_i). \quad (2.23)$$

Note that in this case, since $\check{p}(y_{ij}|u_i)$ is already univariate normal the projection is superfluous. The starting values for $\eta_{p(y_{ij}|u_i) \rightarrow u_i}$ that arises from this substitution is

$$\exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \eta_{p(y_{ij}|u_i) \rightarrow u_i}^{\text{start}} \right\} = \exp \left\{ \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\eta}_{ij} \right\}. \quad (2.24)$$

By matching coefficients of like terms it is easy to show

$$\eta_{p(y_{ij}|u_i) \rightarrow u_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ (2y_{ij} - 1)\zeta'(\hat{a}_{ij}) - \zeta''(\hat{a}_{ij})\hat{u}_i \\ \frac{1}{2}\zeta''(\hat{a}_{ij}) \end{bmatrix} \quad (2.25)$$

where

$$\eta_0^{\text{start}} = \zeta(\hat{a}_{ij}) - (2y_{ij} - 1)\zeta'(\hat{a}_{ij})\hat{u}_i + \frac{1}{2}\zeta''(\hat{a}_{ij})\hat{u}_i^2.$$

Note that in Algorithm [3](#) η_0^{start} is not used in the cycle loop and thus can be set to any arbitrary number without affecting the algorithm. A good choice for \hat{u}_i is Laplace approximation. For the R computing environment, the function `glmer()` of the package “lme4” (Bates, et al., 2018[5](#)) provides fast Laplace approximation-based predictions for the u_i .

2.3 Evaluation of the estimates

Using R software, we implemented and compared the quadrature and EP approach to calculating $\ell(\sigma^2)$ and $\check{\ell}(\sigma^2)$ respectively. Figure [2.2](#) shows estimates of the likelihood surface for both methods. It demonstrates that they provide reasonable estimates of the true value of σ^2 and that any differences between them are extremely small. Having shown equivalence of the methods with regard to estimating the likelihood function, we now turn to methodology for estimation of its maximum with 95% confidence intervals.

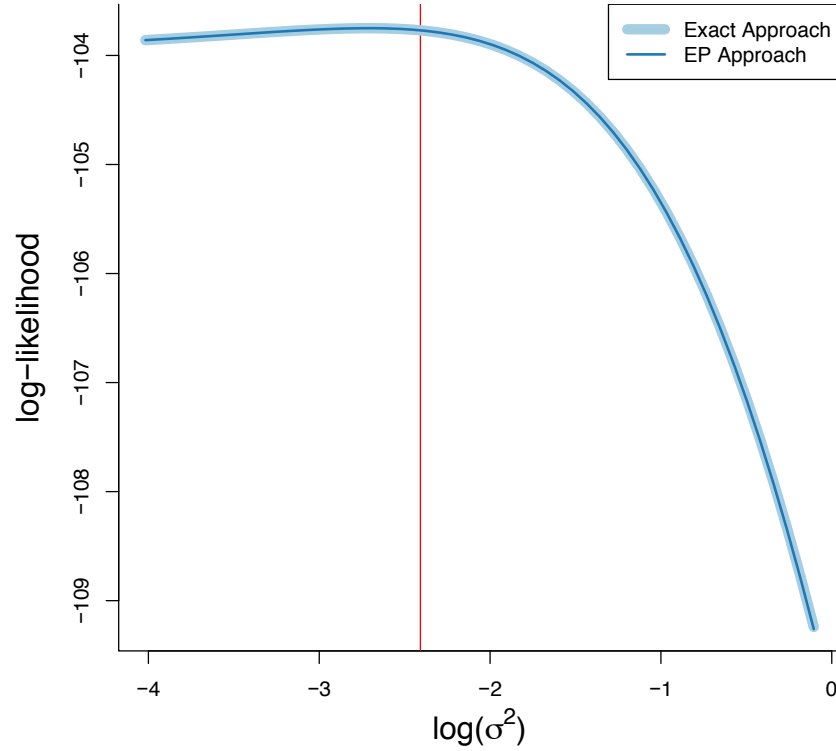


Figure 2.2: A comparison plot of the log-likelihood surface over the parameter σ^2 for probit models calculated using univariate quadrature and EP. The true σ^2 value is represented by the red line on the log scale and other lines follow the legend. The data generated had a 30 groups with 5 responses per group and the true value of σ^2 was 0.09.

2.4 Computing point estimates and confidence intervals

The maximum likelihood estimator of the probit mixed model parameter σ^2 for the quadrature and EP approach are respectively given by $\hat{\sigma}^2 = \operatorname{argmax}_{\sigma^2} \ell(\sigma^2)$ and $\hat{\tilde{\sigma}}^2 = \operatorname{argmax}_{\sigma^2} \tilde{\ell}(\sigma^2)$. To find their stationary points we require the first derivative of the likelihood functions $\ell'(\sigma^2)$ and $\tilde{\ell}'(\sigma^2)$. Calculation of the second derivative (denoted by $\ell''(\sigma^2)$ and $\tilde{\ell}''(\sigma^2)$) facilitates calculation of confidence intervals. Using the second derivatives, we now show how to obtain confidence intervals.

2.4.1 Confidence interval estimation

The overarching schematic for calculation of confidence intervals for both approximate and exact approaches to estimating the σ^2 parameter are analogous so we show it only for the exact case. Since the variance parameter σ^2 is constrained to be a positive number it is more appropriate to work with the parameter in the transformed space

$$\omega \equiv \log(\sigma) = \frac{1}{2} \log(\sigma^2) = g(\sigma^2).$$

Given the maximum likelihood estimator is asymptotically normally distributed, the transformed parameter estimate is

$$\hat{\omega} \sim N \left(g(\sigma_{\text{true}}^2), \frac{1}{\mathbf{I}(g(\hat{\sigma}^2))} \right). \quad (2.26)$$

Since

$$\mathbf{I}(g(\hat{\sigma}^2)) = \frac{\mathbf{I}(\hat{\sigma}^2)}{(g'(\hat{\sigma}^2))^2} \quad (2.27)$$

and $g'(\sigma^2) = (2\sigma^2)^{-1}$, it follows

$$\hat{\omega} \sim N \left(\omega_{\text{true}}, \frac{1}{(2\hat{\sigma}^2)^2 (-\ell''(\hat{\sigma}^2))} \right).$$

Thus for a 95% confidence interval we expect that

$$0.95 \approx P \left(\hat{\omega} - \frac{1.96}{\sqrt{(2\hat{\sigma}^2)^2 (-\ell''(\hat{\sigma}^2))}} < \omega_{\text{true}} < \hat{\omega} + \frac{1.96}{\sqrt{(2\hat{\sigma}^2)^2 (-\ell''(\hat{\sigma}^2))}} \right).$$

Setting

$$\omega_{\text{low}} = \frac{1}{2} \log(\hat{\sigma}^2) - \frac{1.96}{\sqrt{(2\hat{\sigma}^2)^2 (-\ell''(\hat{\sigma}^2))}} \quad \text{and} \quad \omega_{\text{upp}} = \frac{1}{2} \log(\hat{\sigma}^2) + \frac{1.96}{\sqrt{(2\hat{\sigma}^2)^2 (-\ell''(\hat{\sigma}^2))}},$$

the lower and upper 95% confidence intervals for parameter σ^2 are given by

$$\sigma_{\text{low}}^2 = \exp(2\omega_{\text{low}}) \quad \text{and} \quad \sigma_{\text{upp}}^2 = \exp(2\omega_{\text{upp}}).$$

2.4.2 Derivative approximation

For the quadrature and EP approach to obtaining the σ^2 log-likelihood, there are multiple methods by which the first and second derivatives can be obtained. These solutions are either based on analytical or quasi-Newton solutions to the derivatives required.

Quasi-Newton methods negate the need to analytically compute the first and second derivatives and can return the Hessian matrix for higher dimensional cases. We implement both Nelder-Mead (NM) and Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithms for optimisation via the R function `optim()` in the “stats” package (R Core Team, 2019^[56]). We can use unconstrained optimisation on log transformed σ^2 , then back transform to calculate confidence intervals. First we use the more robust NM search to roughly locate the maxima, then refine it via a BFGS search. Using the BFGS method for optimisation at the maxima via the `optim()` function has the additional benefit of returning the second derivative. We initialise the NM algorithm using an estimate for the optimal value $\widehat{\sigma^2}$ via Laplace approximation, which is provided by the R function `glmer()` in the package “lme4” (Bates, et al., 2018)^[5]. Implementation in R of optimisation for both quadrature and EP approaches involves using a simple wrapper function around the parameter to be optimised, such that optimisation is on the unconstrained space.

The analytical solution may be available at a lower computational cost compared to quasi-Newton approaches, however are without doubt more algebraically intensive. First, for the exact likelihood we let

$$J_{si}(\sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \int_{-\infty}^{\infty} u_i^s \prod_{j=1}^n \Phi((2y_{ij} - 1)u_i) \exp(-u_i^2/2\sigma^2) du_i,$$

where $s \in \{0, 1, 2, 3, 4, 5, 6\}$, such that

$$\ell(\sigma^2) = \sum_{i=1}^m \log(J_{0i}(\sigma^2)).$$

Then the first and second derivatives of the likelihood function can be shown to be,

$$\ell'(\sigma^2) = \frac{1}{2\sigma^4} \left(\sum_{i=1}^m \frac{J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} - m\sigma^2 \right) \quad (2.28)$$

and

$$\ell''(\sigma^2) = \frac{1}{4\sigma^8} \left\{ \sum_{i=1}^m \left\{ \frac{J_{4i}(\sigma^2) - 4\sigma^2 J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} - \left(\frac{J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} \right)^2 \right\} + 2m\sigma^4 \right\}. \quad (2.29)$$

In the quadrature case, each $J_{si}(\sigma^2)$ is computed individually before being summed to compute the relevant derivative likelihood. To improve numerical stability we reimplement the work in Section 2.1 to find the maximum value of the integral within each $J_{si}(\sigma^2)$ and limit the range of values each integral can obtain to between 0 and 1, i.e.

$$J_{si}(\sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \int_{-\infty}^{\infty} u_i^s \exp(h(u_i) - h(u_{0i})) du_i \exp(h(u_{0i})), \quad (2.30)$$

where $h(x)$ and u_{0i} are defined as in Section 2.1. Note that the $\exp(h(u_{0i}))$ term becomes superfluous since all $J_{si}(\sigma^2)$ are calculated as ratios of other $J_{si}(\sigma^2)$. We use a bisection search over the surface of $\ell'(\sigma^2)$ to find the minimum and then compute the confidence intervals as explained in Section 2.4.1. Figure 2.3 shows plots of the likelihood and its first and second derivatives obtained via adaptive Gauss-Hermite quadrature.

EP approximation of the maximum likelihood estimate $\widehat{\sigma^2}$ benefits from letting

$$\underline{J}_{si}(\sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \int_{-\infty}^{\infty} u_i^s \psi_i(u_i) \exp(-u_i^2/(2\sigma^2)) du_i,$$

where

$$\psi_i(u) \equiv \exp \left\{ \left[\begin{array}{c} 1 \\ u \\ u^2 \end{array} \right]^\top \text{SUM}(\eta_{p(y_i|u_i) \rightarrow u_i}) \right\}$$

such that

$$\underline{\ell}(\sigma^2) = \sum_{i=1}^m \log(\underline{J}_{0i}(\sigma^2)).$$

Since $\psi_i(u)$ is proportional to a normal density function closed form solutions for the integrals arising in $\underline{\ell}'(\sigma^2)$ and $\underline{\ell}''(\sigma^2)$ exist as before. There are two different approaches for obtaining EP approximations of $\underline{J}_{si}(\sigma^2)$, which we now discuss.

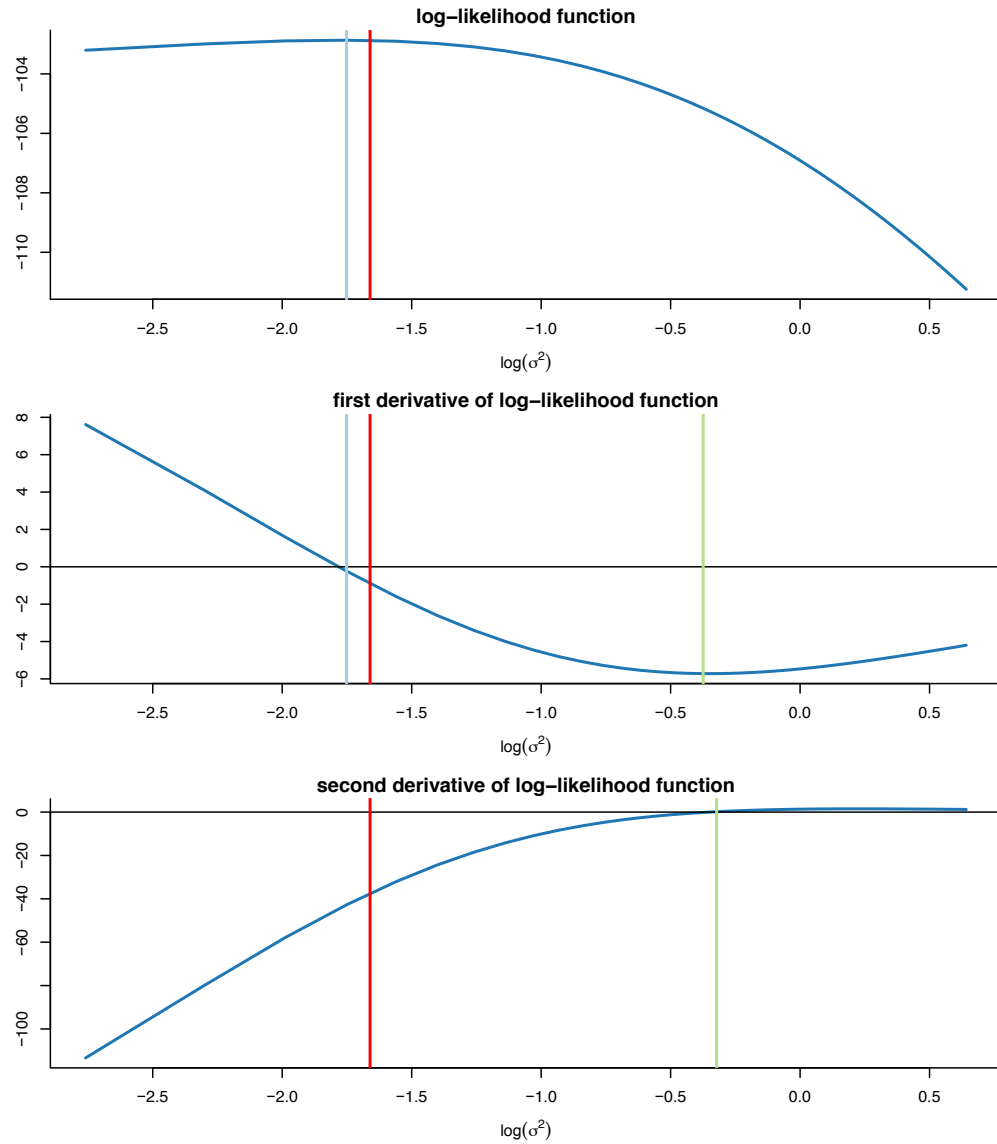


Figure 2.3: Plot of the $\ell(\sigma^2)$ and its first two derivatives solved via adaptive Gauss-Hermite quadrature. The red line represents the true $\sigma^2 = 0.19$, the light blue line is the point where the first derivative cuts the x axis and the green line is analogous to the light blue line for the second derivative. The number of groups was 30 with 5 responses per group.

2.4.2.1 Expectation propagation analytical approach I

The first approach involves directly approximating the $J_{si}(\sigma^2)$ via the natural parameters passed in Algorithm 3. Consider that,

$$J_{si}(\sigma^2) = \int_{-\infty}^{\infty} u^s \exp \left\{ \left[\begin{array}{c} 1 \\ u \\ u^2 \end{array} \right]^{\top} \boldsymbol{\eta}_i^{\boxplus} \right\} du,$$

where

$$\boldsymbol{\eta}_i^{\boxplus} = \begin{bmatrix} \eta_{0i}^{\boxplus} \\ \eta_{1i}^{\boxplus} \\ \eta_{2i}^{\boxplus} \end{bmatrix} \equiv \boldsymbol{\eta}_{\sigma^2} + \text{SUM}(\boldsymbol{\eta}_{p(y_i|u_i) \rightarrow u_i}) \quad (2.31)$$

and $\boldsymbol{\eta}_{\sigma^2}$ is defined in equation (2.8). Since u_i is a normal random variable with natural parameters η_{1i}^{\boxplus} and η_{2i}^{\boxplus} , we can write

$$J_{si}(\sigma^2) = (2\pi)^{-1/2} \exp(\eta_{0i}^{\boxplus} u + A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})) E(u^s),$$

where

$$E(u^s) = (2\pi)^{-1/2} \int_{-\infty}^{\infty} u^s \exp(\eta_{1i}^{\boxplus} u + \eta_{2i}^{\boxplus} u^2 - A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})) du.$$

The integrals required for these expressions can be solved with closed form solutions. Consider Result 7.

Result 7. If x is a normal random variable with natural parameters η_1 and η_2 , then

$$\begin{aligned} E(x) &= -\frac{\eta_1}{2\eta_2}, \\ E(x^2) &= \frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2\eta_2}, \\ E(x^3) &= \frac{3\eta_1}{4\eta_2^2} - \frac{\eta_1^3}{8\eta_2^3}, \\ E(x^4) &= \frac{\eta_1^4}{16\eta_2^4} - \frac{3\eta_1^2}{4\eta_2^3} + \frac{3}{4\eta_2^2}, \\ E(x^5) &= \frac{5\eta_1^3}{8\eta_2^4} - \frac{\eta_1^5}{32\eta_2^5} - \frac{15\eta_1}{8\eta_2^3}, \end{aligned}$$

and

$$E(x^6) = \frac{\eta_1^6}{64\eta_2^6} - \frac{15\eta_1^4}{32\eta_2^5} + \frac{45\eta_1^2}{16\eta_2^4} - \frac{15}{8\eta_2^3}.$$

By Result [7](#),

$$\begin{aligned} \mathcal{J}_{0i}(\sigma^2) &= (2\pi)^{-1/2} \exp(\eta_{0i}^{\boxplus} + A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})), \\ \mathcal{J}_{1i}(\sigma^2) &= (2\pi)^{-1/2} \exp(\eta_{0i}^{\boxplus} + A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})) (-\eta_{1i}^{\boxplus}/(2\eta_{2i}^{\boxplus})), \\ \mathcal{J}_{2i}(\sigma^2) &= (2\pi)^{-1/2} \exp(\eta_{0i}^{\boxplus} + A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})) ((\eta_{1i}^{\boxplus})^2 - 2\eta_{2i}^{\boxplus})/(4(\eta_{2i}^{\boxplus})^2) \end{aligned}$$

and

$$\mathcal{J}_{4i}(\sigma^2) = (2\pi)^{-1/2} \exp(\eta_{0i}^{\boxplus} + A(\eta_{1i}^{\boxplus}, \eta_{2i}^{\boxplus})) \left\{ (\eta_{1i}^{\boxplus})^4 + 12\eta_{2i}^{\boxplus}(\eta_{2i}^{\boxplus} - (\eta_{1i}^{\boxplus})^2) \right\} / (16(\eta_{2i}^{\boxplus})^4).$$

To aid with simplification of the following calculations we introduce Result [8](#).

Result 8. When $s = \{1, 2, 4\}$, ratios of the form $\frac{\mathcal{J}_{si}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)}$ can be simplified to

$$\begin{aligned} \frac{\mathcal{J}_{1i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} &= -\eta_{1i}^{\boxplus}/(2\eta_{2i}^{\boxplus}), \\ \frac{\mathcal{J}_{2i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} &= ((\eta_{1i}^{\boxplus})^2 - 2\eta_{2i}^{\boxplus})/(4(\eta_{2i}^{\boxplus})^2), \\ \frac{\mathcal{J}_{4i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} &= \left\{ (\eta_{1i}^{\boxplus})^4 + 12\eta_{2i}^{\boxplus}(\eta_{2i}^{\boxplus} - (\eta_{1i}^{\boxplus})^2) \right\} / (16(\eta_{2i}^{\boxplus})^4). \end{aligned}$$

Recalling the form of the exact expression from equation (2.28), we can implement an EP approximation approach to the first derivative by replacing each J_{si} with \underline{J}_{si} . This leads to

$$\underline{\ell}'(\sigma^2) = \frac{1}{2\sigma^4} \left(\sum_{i=1}^m \frac{\underline{J}_{2i}(\sigma^2)}{\underline{J}_{0i}(\sigma^2)} - m\sigma^2 \right).$$

Analogously, the EP approximation approach to the second derivative can be obtained by replacing each J_{si} by \underline{J}_{si} in the $\ell''(\sigma^2)$ expression given in equation (2.29), such that

$$\underline{\ell}''(\sigma^2) = \frac{1}{4\sigma^8} \left\{ \sum_{i=1}^m \left\{ \frac{J_{4i}(\sigma^2) - 4\sigma^2 J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} - \left(\frac{J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} \right)^2 \right\} + 2m\sigma^4 \right\}.$$

In summary, at the end of Algorithm 2.2.3 when the natural parameters have converged, they can be used to calculate the derivatives required. Algorithm 4 aims to clarify this process.

Algorithm 4 Alteration of Algorithm 3 to calculate the log-likelihood of σ^2 and its first and second derivatives via EP analytical approach I.

After convergence is reached of Algorithm 3, we can obtain $\underline{\ell}'(\sigma^2)$ and $\underline{\ell}''(\sigma^2)$ by,

$$\begin{aligned} \underline{\ell}'(\sigma^2) &= \frac{1}{2\sigma^4} \left(\sum_{i=1}^m \frac{\underline{J}_{2i}(\sigma^2)}{\underline{J}_{0i}(\sigma^2)} - m\sigma^2 \right), \\ \underline{\ell}''(\sigma^2) &= \frac{1}{4\sigma^8} \left\{ \sum_{i=1}^m \left\{ \frac{J_{4i}(\sigma^2) - 4\sigma^2 J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} - \left(\frac{J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} \right)^2 \right\} + 2m\sigma^4 \right\} \end{aligned}$$

where the ratios of \underline{J}_{si} are defined in Result 8.

2.4.2.2 Expectation propagation analytical approach II

The second approach involves approximating each $\ell'_i(\sigma^2)$ and $\ell''_i(\sigma^2)$ by finding the required $J_{si}(\sigma^2)$ and then summing them to obtain the first and second derivatives of the likelihood. As opposed to the previous approach where we project $f_{\text{input}}(x)$ (as defined in equation (2.5)) onto the normal family, here we project it onto what we refer to as the ‘‘power normal’’ family, which has the general form of

$$p(x) \propto x^s (2\pi\sigma^2)^{-1/2} \exp \left(- (x - \mu)^2 / (2\sigma^2) \right), \quad s \in \{0, 2, 4\}.$$

As before, closed form solutions to the integrals required exist, however the algebra necessary becomes increasingly difficult as s increases. Consider then, the following family of unnormalised power normal density functions written in exponential family form

$$f_{\text{UN}}(x) = \exp \left\{ \begin{bmatrix} x^s \\ x^{s+1} \\ x^{s+2} \end{bmatrix}^\top \begin{bmatrix} \eta_0 \\ \eta_1 \\ \eta_2 \end{bmatrix} \right\}, \quad (2.32)$$

with natural parameters $\eta_0, \eta_1 \in \mathbb{R}$ and $\eta_2 < 0$. Then as before, we wish to find

$$\text{proj}_{\text{UN}}[f_{\text{input}}](x) = \exp \left\{ \begin{bmatrix} x^s \\ x^{s+1} \\ x^{s+2} \end{bmatrix}^\top \boldsymbol{\eta}^* \right\}, \quad (2.33)$$

where

$$\boldsymbol{\eta}^* \equiv \begin{bmatrix} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{bmatrix}$$

and

$$(\eta_0^*, \eta_1^*, \eta_2^*) = \underset{(\eta_0, \eta_1, \eta_2) \in H}{\text{argmin}} \text{KL}(f_{\text{input}} \parallel f_{\text{UN}}).$$

For the first and second derivatives of probit binary GLMMs, EP requires repeated projection of the form given in equation (2.5) onto an unnormalised power normal distribution where $\eta_1^{\text{input}} \in \mathbb{R}$ and $\eta_2^{\text{input}} < 0$. Consider Result 9.

Result 9. For an unnormalised input function $f \in L_1$ such that $f(x) \geq 0$ for all $x \in \mathbb{R}$ where $C_f \equiv \int_{\mathbb{R}} f(x) dx$, the projection onto the unnormalised power normal family is

$$\text{proj}_{\text{UPN}}[f] = C_f \text{proj}_{\text{PN}}[f/C_f].$$

where $\text{proj}_{\text{PN}}[\cdot]$ is the projection onto the power normal family and $\text{proj}_{\text{UPN}}[\cdot]$ is the projection onto the unnormalised power normal family.

As before, the optimal natural parameters η_1^* and η_2^* are given according to the projection of the normalised function f/C_f onto the power normal family. We can subsequently use these optimal natural parameters to find the normalising natural

parameter η_0^* via Result 4 and thus obtain the projection onto unnormalised power normal family.

Result 10. *When the input density follows the form of equation (2.5), projections onto the power normal family η_0^* is given by*

$$\eta_0^* = \log(C_f) - A_s(\eta_1^*, \eta_2^*) + \log h_s(x),$$

where

$$A_s(\eta_1, \eta_2) = \frac{-\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) + \log(E(x^s)), \quad h_s(x) = \frac{x^s}{\sqrt{2\pi}}$$

and $E(x^s)$ is given by Result 7.

Thus to obtain the required projection, we first obtain the optimal natural parameters η_1^* and η_2^* to project onto the power normal family. To do this, first consider the normalised approximation of the input density to be of the form

$$f_{\text{input}}(x) = x^s (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) (E(x^s))^{-1},$$

where

$$E(x^s) = \int_{-\infty}^{\infty} x^s (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx$$

is used to normalise the unnormalised approximation of $f_{\text{input}}(x)$ from the power normal family. Next note that using matrix notation and converting to natural parameters,

$$f_{\text{input}}(x) = x^s \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} - \left\{ -\frac{1}{2} \log(-2\eta_2) + \log(E(x^s)) - \frac{\eta_1^2}{4\eta_2} \right\} \right\} (2\pi)^{-1/2}.$$

As in the proof of Result 5, this can be expressed as

$$f_{\text{input}}(x) = \exp(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1} - A_s(\boldsymbol{\eta}_{-1})) h_s(x),$$

where

$$\mathbf{T}(x) \equiv \begin{bmatrix} x \\ x^2 \end{bmatrix}, \quad h_s(x) \equiv \frac{x^s}{\sqrt{2\pi}}, \quad \boldsymbol{\eta}_{-1} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \equiv \begin{bmatrix} \mu/\sigma^2 \\ -1/(2\sigma^2) \end{bmatrix} \quad (2.34)$$

and

$$A_s(\boldsymbol{\eta}_{-1}) = \frac{-\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) + \log\{E(x^s)\}.$$

Consider the following relationship which defines the gradient of the likelihood function at $\boldsymbol{\eta}_{-1}$:

$$\nabla f(\boldsymbol{\eta}_{-1}) \equiv (\nabla A_s)(\boldsymbol{\eta}_{-1}) - \boldsymbol{\tau}. \quad (2.35)$$

Note that when $\nabla f(\boldsymbol{\eta}_{-1}) = 0$ we have

$$\boldsymbol{\eta}_{-1} \equiv (\nabla A_s)^{-1}(\boldsymbol{\tau})$$

where

$$\boldsymbol{\tau} = \begin{bmatrix} \mathcal{M}_{s+1}/\mathcal{M}_s \\ \mathcal{M}_{s+2}/\mathcal{M}_s \end{bmatrix}, \quad (\nabla A_s)(\boldsymbol{\eta}_{-1}) = \begin{bmatrix} \frac{\partial}{\partial \eta_1} A_s(\boldsymbol{\eta}_{-1}) \\ \frac{\partial}{\partial \eta_2} A_s(\boldsymbol{\eta}_{-1}) \end{bmatrix}$$

and \mathcal{M}_s is defined as per Section [2.7.2](#) of the appendix. To find the optimal natural parameters for the projection required we must minimise the function $(\nabla A_s)(\boldsymbol{\eta}) - \boldsymbol{\tau}$ such that we obtain the $\boldsymbol{\eta}_{-1}^*$ where it is 0. We do this using a Newton-Raphson search, initialised at a rough estimate of the value that minimises the sum of squares deviation from the root of the inverse-defining equation,

$$\boldsymbol{\eta}_{-1}^* = \operatorname{argmin} \sum ((\nabla A_s)(\boldsymbol{\eta}_{-1}) - \boldsymbol{\tau})^2,$$

found via a Nelder-Mead search initialised at $(\nabla A_s)(\boldsymbol{\eta})$ for the normal family. We can accomplish the Nelder-Mead step using the same `optim()` function in R as in the previous section. To match our previous notation, consider Definition [12](#).

Definition 12. Define the function

$$k_{\text{probitPower}}(\mathbf{a}; c_0, c_1, s) = (\nabla A)^{-1}(\boldsymbol{\tau}, s), \quad (2.36)$$

where

$$\boldsymbol{\tau} = \begin{bmatrix} \mathcal{M}_{s+1}(c_0, c_1, \mathbf{a}) / \mathcal{M}_s(c_0, c_1, \mathbf{a}) \\ \mathcal{M}_{s+2}(c_0, c_1, \mathbf{a}) / \mathcal{M}_s(c_0, c_1, \mathbf{a}) \end{bmatrix},$$

$(\nabla A)^{-1}(\boldsymbol{\tau}, s)$ is the function that returns $\boldsymbol{\eta}_{-1}$ when $\nabla f(\boldsymbol{\eta}_{-1}) = 0$ (as per equation (2.35)) and $\mathcal{M}_{s+k}(c_0, c_1, \mathbf{a})$ is

$$\mathcal{M}_0(c_0, c_1, \mathbf{a}) = \mathcal{W}_0(r_8, r_9),$$

$$\mathcal{M}_1(c_0, c_1, \mathbf{a}) = a_1 \mathcal{W}_0(r_8, r_9) + c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9),$$

$$\mathcal{M}_2(c_0, c_1, \mathbf{a}) = a_1^2 \mathcal{W}_0(r_8, r_9) + 2a_1 c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9) - 2a_2 \mathcal{W}_2(r_8, r_9),$$

$$\mathcal{M}_3(c_0, c_1, \mathbf{a}) = a_1^3 \mathcal{W}_0(r_8, r_9) + 3a_1^2 c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9) - 6a_1 a_2 \mathcal{W}_2(r_8, r_9) - 2c_1 r_9^{-1} a_2 \mathcal{W}_3(r_8, r_9),$$

$$\mathcal{M}_4(c_0, c_1, \mathbf{a}) = a_1^4 \mathcal{W}_0(r_8, r_9) + 4a_1^3 c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9) - 12a_1^2 a_2 \mathcal{W}_2(r_8, r_9) - 8a_1 a_2 c_1 r_9^{-1} \mathcal{W}_3(r_8, r_9) + 4a_2^2 \mathcal{W}_4(r_8, r_9),$$

$$\mathcal{M}_5(c_0, c_1, \mathbf{a}) = a_1^5 \mathcal{W}_0(r_8, r_9) + 5a_1^4 c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9) - 20a_1^3 a_2 \mathcal{W}_2(r_8, r_9) - 20a_1^2 a_2 c_1 r_9^{-1} \mathcal{W}_3(r_8, r_9) + 20a_1 a_2^2 \mathcal{W}_4(r_8, r_9) + 4a_2^2 c_1 r_9^{-1} \mathcal{W}_5(r_8, r_9),$$

$$\mathcal{M}_6(c_0, c_1, \mathbf{a}) = a_1^6 \mathcal{W}_0(r_8, r_9) + 6a_1^5 c_1 r_9^{-1} \mathcal{W}_1(r_8, r_9) - 30a_1^4 a_2 \mathcal{W}_2(r_8, r_9) - 40a_1^3 a_2 c_1 r_9^{-1} \mathcal{W}_3(r_8, r_9) + 60a_1^2 a_2^2 \mathcal{W}_4(r_8, r_9) + 24a_1 a_2^2 c_1 r_9^{-1} \mathcal{W}_5(r_8, r_9) - 8a_2^3 \mathcal{W}_6(r_8, r_9),$$

for $\mathcal{W}_{s+k}(a, b)$ defined as in equation (2.47) and (2.43), $\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$, $r_1 =$

$$\sqrt{2(2 - c_1^2 a_2^{-1})}, r_2 = (2c_0 - c_1 a_2^{-1} a_1) r_1^{-1}, r_8 = 2r_2 r_1 \text{ and } r_9 = c_1 (-2a_2)^{-1/2}.$$

Note that the derivation of Definition 12 follows that of Definition 5 and as such we do not provide details and instead refer interested readers to Appendix 2.7.2. Each $\mathcal{M}_s(c_0, c_1, \mathbf{a})$ is found by simple algebra in an analogous manner to those required for Definition 10. The first and second derivatives required for both search methods of $(\nabla A)^{-1}(\boldsymbol{\tau}, s)$ are provided in Appendix 2.7.4. As before, to obtain the required projection, we present Result 11.

Result 11. Given f follows the form of equation (2.5), the projection onto the power normal family is given by

$$proj_{PN}[f_{input}] = \exp(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1}^* - A_s(\boldsymbol{\eta}_{-1}^*)) h_s(x)$$

where

$$\boldsymbol{\eta}_{-1}^* = k_{probitPower}(\boldsymbol{\eta}^{input}; c_0, c_1, s), \quad \boldsymbol{\eta}_{-1}^{input} \equiv \begin{bmatrix} \eta_1^{input} \\ \eta_2^{input} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix},$$

and $\mathbf{T}(x)$ and $h_s(x)$ follow from equation (2.34).

We now require a function to find the normalising natural parameter η_0^* required for the projection onto unnormalised power normal family. Following analogous arguments used for the projection onto the unnormalised normal family, we define

$$c_{probitPower} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}; c_0, c_1, s \right) = \log \mathcal{M}_s \left(c_0, c_1, \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \right) - A_s \left(\begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \right) - \log h_s(x),$$

where $\log(\mathcal{M}_s(c_0, c_1, \mathbf{a}))$, $A_s(\mathbf{b})$ and $\log h_s(x)$ are defined as per the previous definitions. We present Algorithm 5 which requires three loops over the main portion of the algorithm to calculate the each $J_{si}(\sigma^2)$, as opposed to the previous method. Additionally note that each $J_{si}(\sigma^2)$ is found in each loop and then used to find the first and second derivative of the likelihood.

Algorithm 5 Algorithm to calculate the log-likelihood of σ^2 and its first and second derivatives via EP analytical approach II.

Set constants: $c_0 \leftarrow 0$; $c_{1ij} \leftarrow 2y_{ij} - 1$; $1 \leq i \leq m$, $1 \leq j \leq n$.

$$\boldsymbol{\eta}_{p(u_i; \sigma^2) \rightarrow u_i} \leftarrow \boldsymbol{\eta}_{\sigma^2} \equiv \begin{bmatrix} -\frac{1}{2} \log(2\pi\sigma^2) \\ 0 \\ 1/(2\sigma^2) \end{bmatrix}, \quad 1 \leq i \leq m.$$

For $s = 0, 2, 4$:

For $i = 1, \dots, m$:

Initialise: $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \in \mathbb{R}$, $1 \leq j \leq n$ as per the equation (2.25).

Cycle:

$$\text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) \leftarrow \sum_{j=1}^n (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})$$

For $j = 1, \dots, n$:

$$\begin{aligned} \boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)} &\leftarrow \boldsymbol{\eta}_{\sigma^2} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}) - \boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i} \\ (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} &\leftarrow k_{\text{probitPower}}\left((\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}; c_0, c_{1ij}, s\right) \\ &\quad - (\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2} \end{aligned}$$

until all natural parameter vectors converge.

For $j = 1, \dots, n$:

$$\begin{aligned} (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 &\leftarrow c_{\text{probitPower}}\left((\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}, (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2}\right. \\ &\quad \left. + (\boldsymbol{\eta}_{u_i \rightarrow p(y_{ij}|u_i)})_{1:2}; c_0, c_{1ij}, s\right) \end{aligned}$$

$$\boldsymbol{\eta}_i^{\boxplus} \equiv \boldsymbol{\eta}_{\sigma^2} + \text{SUM}(\boldsymbol{\eta}_{p(y_i|u_i) \rightarrow u_i})$$

$$\mathcal{J}_{0i}(\sigma^2) = (2\pi)^{-1/2} \exp(\boldsymbol{\eta}_{0i}^{\boxplus} + A(\boldsymbol{\eta}_{1i}^{\boxplus}, \boldsymbol{\eta}_{2i}^{\boxplus})),$$

$$\text{If } s = 2, \quad \mathcal{J}_{2i}(\sigma^2) = \mathcal{J}_{0i}(\sigma^2) \left((\boldsymbol{\eta}_{1i}^{\boxplus})^2 - 2\boldsymbol{\eta}_{2i}^{\boxplus} \right) / (4(\boldsymbol{\eta}_{2i}^{\boxplus})^2),$$

$$\text{If } s = 4, \quad \mathcal{J}_{4i}(\sigma^2) = \mathcal{J}_{0i}(\sigma^2) \left\{ (\boldsymbol{\eta}_{1i}^{\boxplus})^4 + 12\boldsymbol{\eta}_{2i}^{\boxplus}(\boldsymbol{\eta}_{2i}^{\boxplus} - (\boldsymbol{\eta}_{1i}^{\boxplus})^2) \right\} / (16(\boldsymbol{\eta}_{2i}^{\boxplus})^4).$$

where the ratios of \mathcal{J}_{si} are defined in Result 8.

After convergence is reached.

$$\underline{\ell}'(\sigma^2) = \frac{1}{2\sigma^4} \left(\sum_{i=1}^m \frac{\mathcal{J}_{2i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} - m\sigma^2 \right),$$

$$\underline{\ell}''(\sigma^2) = \frac{1}{4\sigma^8} \left\{ \sum_{i=1}^m \left\{ \frac{\mathcal{J}_{4i}(\sigma^2) - 4\sigma^2 \mathcal{J}_{2i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} - \left(\frac{\mathcal{J}_{2i}(\sigma^2)}{\mathcal{J}_{0i}(\sigma^2)} \right)^2 \right\} + 2m\sigma^4 \right\}.$$

2.5 Best predictor

In addition to estimating σ^2 , we also wish to predict the random effect u_i for each group. First, consider that in the exact case

$$\begin{aligned} \text{BP}(u_i) &= E(u_i|\mathbf{y}_i) \\ &= \int_{-\infty}^{\infty} u_i p(u_i|\mathbf{y}_i; \sigma^2) du_i \\ &= \frac{\int_{-\infty}^{\infty} u_i p(\mathbf{y}_i|u_i) p(u_i; \sigma^2) du_i}{\int_{-\infty}^{\infty} p(\mathbf{y}_i|u_i) p(u_i; \sigma^2) du_i} = \frac{J_{1i}(\sigma^2)}{J_{0i}(\sigma^2)}, \end{aligned} \quad (2.37)$$

where $\mathbf{y}_i = (y_{i1}, \dots, y_{in})$ and $J_{si}(\sigma^2)$ for $s \in \{0, 1, 2\}$ follows from equation (2.28). Equation (9.7) of McCulloch, Searle & Neuhaus (2008)^[41] suggests that $E_{\mathbf{y}_i}(\text{Var}(u_i|\mathbf{y}_i))$ provides a reasonable approximation of prediction errors i.e. $\text{Var}(\text{BP}(u_i) - u_i)$, where

$$\text{Var}(u_i|\mathbf{y}_i) = \frac{J_{2i}(\sigma^2)}{J_{0i}(\sigma^2)} - \left(\frac{J_{1i}(\sigma^2)}{J_{0i}(\sigma^2)} \right)^2.$$

However, the approximation is hindered by the expectation over the distribution of the \mathbf{y}_i vector. For the EP case, the best predictor of u_i can be obtained via products of Algorithm [3](#). Let

$$\hat{\boldsymbol{\eta}}_i \equiv \boldsymbol{\eta}_\sigma^2 + \text{SUM}(\boldsymbol{\eta}_{p(\mathbf{y}_i|u_i) \rightarrow u_i}) = \begin{bmatrix} \hat{\eta}_{i1} \\ \hat{\boldsymbol{\eta}}_{i2} \end{bmatrix}, \quad (2.38)$$

where $\hat{\eta}_{i1}$ corresponds to the first entry of $\hat{\boldsymbol{\eta}}_i$, $\hat{\boldsymbol{\eta}}_{i2}$ corresponds to the remaining entries and $\boldsymbol{\eta}_\sigma^2$ and $\text{SUM}(\boldsymbol{\eta}_{p(\mathbf{y}_i|u_i; \boldsymbol{\beta}) \rightarrow u_i})$ are as previously defined in Algorithm [3](#) with $\sigma^2 = \hat{\sigma}^2$. Note that Algorithm [3](#) involves using

$$\exp \left\{ \begin{bmatrix} 1 \\ u_i \\ u_i^2 \end{bmatrix}^\top \hat{\boldsymbol{\eta}}_i \right\} \quad \text{to replace} \quad p(\mathbf{y}_i|u_i) p(u_i; \sigma^2).$$

We can approximate $\text{BP}(u_i) = E(\hat{u}_i)$, where \hat{u}_i is univariate normal with natural parameters $\hat{\boldsymbol{\eta}}_i$. Thus,

$$\text{BP}(u_i) = -\hat{\eta}_{i1} / 2\hat{\boldsymbol{\eta}}_{i2}.$$

As before, the approximation of prediction errors given by $E_{\mathbf{y}_i}(\text{Var}(u_i|\mathbf{y}_i))$, where

$$\text{Var}(u_i|\mathbf{y}) = -1/2\hat{\boldsymbol{\eta}}_{i2},$$

is hindered by the expectation over the distribution over the \mathbf{y}_i vector.

2.6 Simulation Study

A simulation study comparing four methods of obtaining maximum likelihood estimates and confidence intervals in Section 2.4 was run in the R computing environment (R Core Team, 2019^[56]) on a MacBook Air laptop with two 2.2 gigahertz processors and 8 gigabytes of random access memory. The four methods are:

- Adaptive Gauss-Hermite quadrature using quasi-Newton optimisation and second derivative calculation.
- EP using quasi-Newton optimisation and second derivative calculation.
- EP using analytical approach I to first and second derivative calculation.
- EP using analytical approach II to first and second derivative calculation.

Preliminary studies of datasets of $m = 20$ groups with $n = 2$ observations per group showed a lack of stability for EP analytical approach II. It was deemed such instability was caused by poor starting values in the bisection search. Although it is possible to increase the number of iterations in the bisection and get convergence on a root, for speed starting values were based on the results of the quadrature approach using quasi-Newton methods. Additionally, poor behaviour of the likelihood function for low sample sizes may have caused problems during the optimisation of its derivative.

With this in mind our study consisted of 100 repetitions of $\hat{\sigma}$ estimates for each of the estimation methods for each set of simulated data. The true σ value of the dataset was either $\sigma_{\text{true}} = 0.25, 0.50$ or 1.00 . Each of these true values were used to generate test datasets with $m = 50, 250$ and 1250 groups. The number of observations per group was fixed at $n = 10$. Due to the presence of multiple iterative loops, computational speed was severely compromised. Since the code used implements a variety of R packages which make calls to various low level languages, we omit a comparison of run speed and focus purely on the accuracy and stability of the methods.

A random sample of 10 confidence intervals for each method are compared visually in Figure 2.5, showing although there is generally parity between the methods, approach I to EP returned the tightest intervals by a small margin.

A comparison between the estimated $\hat{\sigma}$ of each method and the σ_{true} for each dataset was conducted using Wilcoxon-tests. The p-values from this study are presented with

boxplots of log estimate error in Figure 2.4, where the error ϵ is given by

$$\epsilon = \frac{\|\hat{\sigma} - \sigma_{\text{true}}\|}{\hat{\sigma}}. \quad (2.39)$$

For each of the methods tested and each of the σ_{true} values, the variance decreased as the number of groups in the dataset increased. For the dataset where $\sigma_{\text{true}} = 0.5$, there was no statistical difference between the estimate and true parameter regardless of the number of groups in the dataset. For the dataset with $\sigma_{\text{true}} = 1$ and $m = 1250$ groups, the variance of the estimates for the exact case using `optim()` and approach I to EP increased, while the other methods did not. For datasets where $\sigma_{\text{true}} = 0.25$, there was no statistical difference for group sizes less than $m = 1250$. The discrepancy at $m = 1250$ is most likely due to the reduction of variance in datasets with a high number of groups which lead to an oversensitive test statistic.

Given the algebraic simplifications afforded by using quasi-Newton methods to facilitate derivative calculation in addition to the stability problems of other methods, there is little evidence it is worthwhile to continue with Approach I and Approach II to EP. Furthermore these approaches are difficult to implement for the general case of GLMMs where variates and random effects can be any dimension.

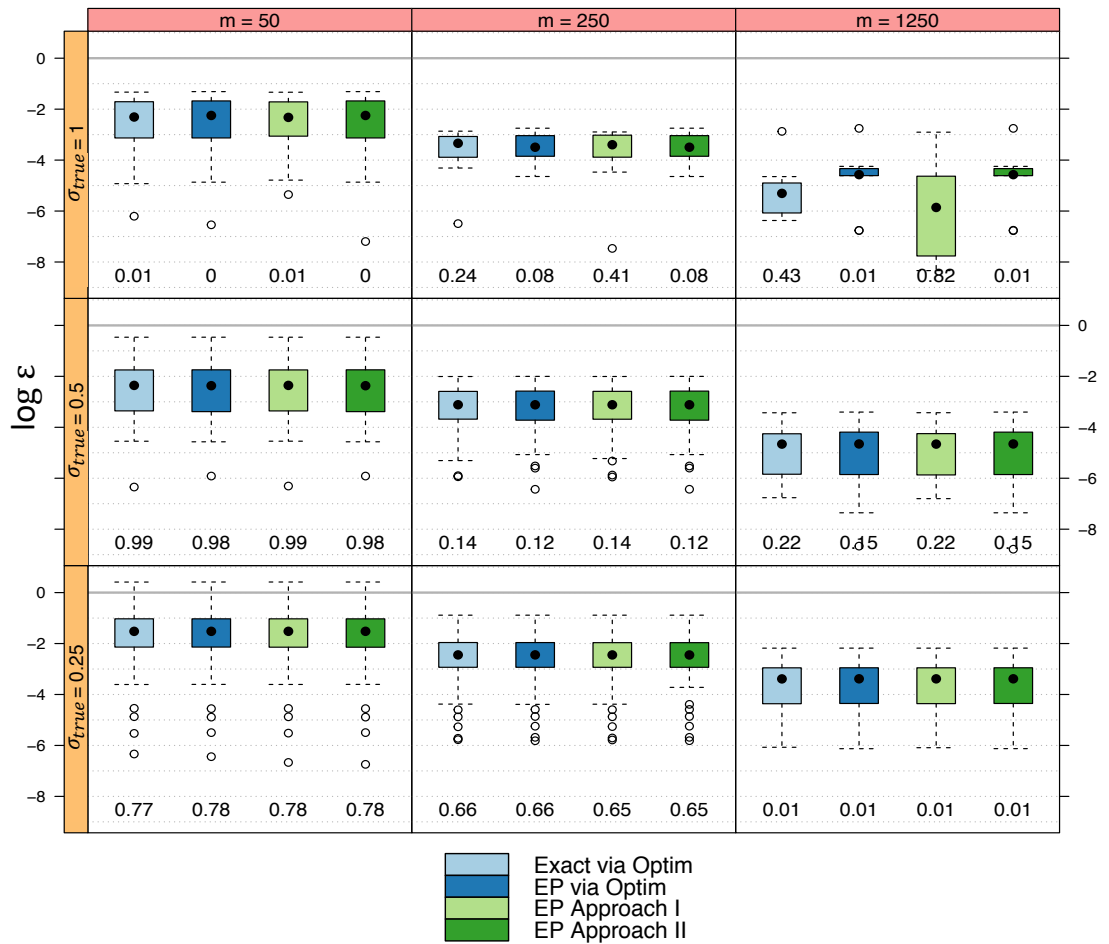


Figure 2.4: Panel plot with error boxplots for the four methods of approximation discussed in this section. P-values from a Wilcoxon comparison test are shown under each boxplot. The legend shows which method produced each boxplot. The calculation of error is given in equation (2.39).

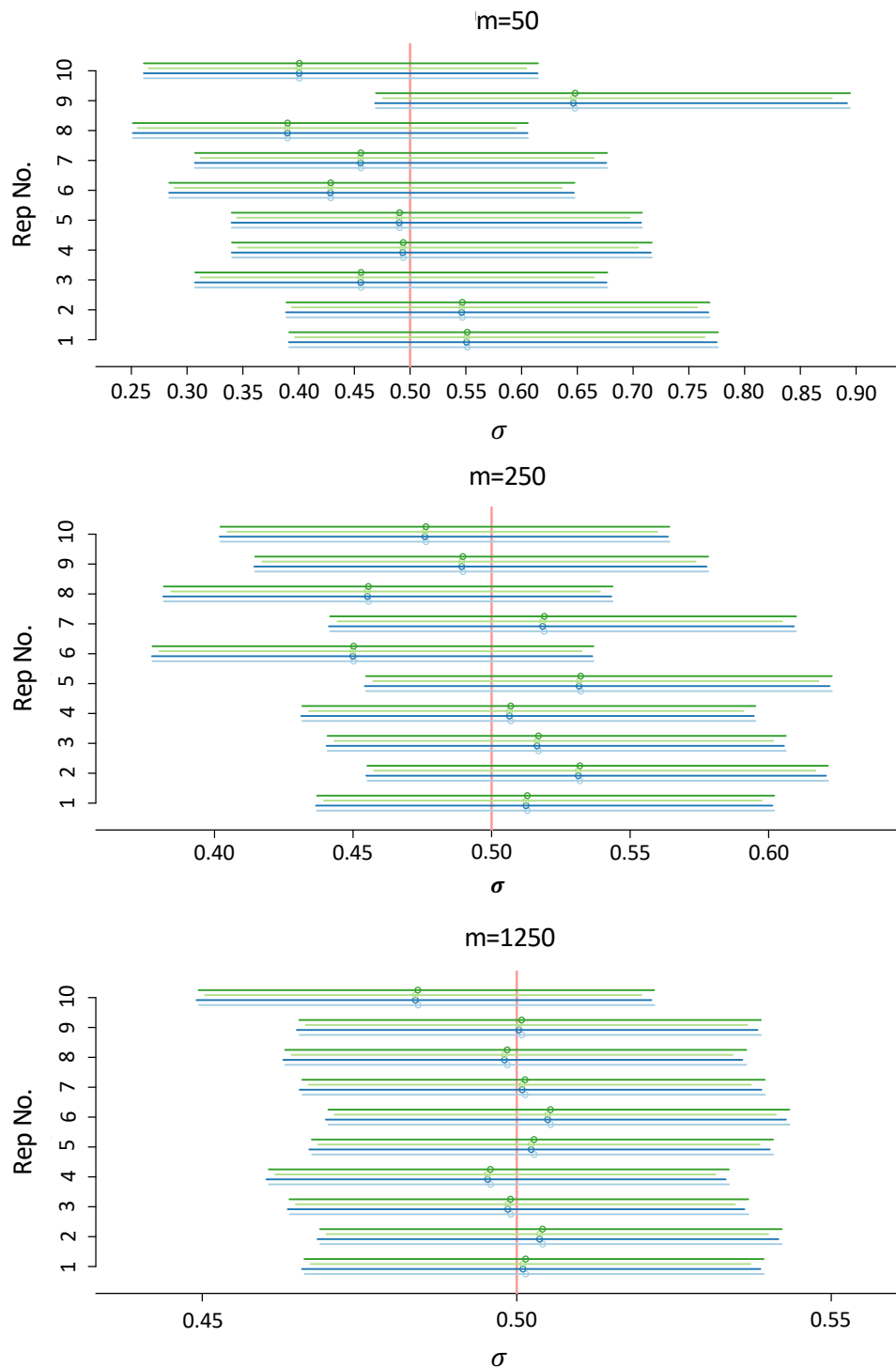


Figure 2.5: Plot displaying 10 randomly selected 95% confidence intervals of each method across various group sizes and $\sigma_{true} = 0.5$, as shown by the red line. The lines correspond with the methods of this section as shown in the legend of Figure [2.4](#).

2.7 Appendix

2.7.1 Proof of Result 3

Consider an unnormalised normal univariate density function

$$f_{\text{UN}}(x; \boldsymbol{\eta}) = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta} \right\}.$$

Then the KL divergence of $f_{\text{UN}}(x; \boldsymbol{\eta})$ from f_{input} is

$$\begin{aligned} \text{KL}(f_{\text{input}} \parallel f_{\text{UN}}) &= \int_{\mathbb{R}} f_{\text{input}}(x) \log(f_{\text{input}}(x)/f_{\text{UN}}(x; \boldsymbol{\eta})) + f_{\text{UN}}(x; \boldsymbol{\eta}) - f_{\text{input}}(x) dx \\ &= \mathcal{K}(\boldsymbol{\eta}) + \text{const}, \end{aligned}$$

where

$$\mathcal{K}(\boldsymbol{\eta}) \equiv (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A(\boldsymbol{\eta}_{1:2})) - \begin{bmatrix} \int_{\mathbb{R}} f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x^2 f_{\text{input}}(x) dx \end{bmatrix}^\top \boldsymbol{\eta},$$

and “const” represents all the terms not dependent on $\boldsymbol{\eta}$. Note that the derivative vector of $\mathcal{K}(\boldsymbol{\eta})$ is

$$D\mathcal{K}(\boldsymbol{\eta}) = (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A(\boldsymbol{\eta}_{1:2})) \begin{bmatrix} 1 \\ DA(\boldsymbol{\eta}_{1:2})^\top \end{bmatrix}^\top - \begin{bmatrix} \int_{\mathbb{R}} f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x^2 f_{\text{input}}(x) dx \end{bmatrix}^\top.$$

Setting $D\mathcal{K}(\boldsymbol{\eta})^\top = 0$ to minimise the $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$,

$$(2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A(\boldsymbol{\eta}_{1:2})) \begin{bmatrix} 1 \\ \nabla A(\boldsymbol{\eta}_{1:2}) \end{bmatrix}^\top = \begin{bmatrix} \int_{\mathbb{R}} f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x f_{\text{input}}(x) dx \\ \int_{\mathbb{R}} x^2 f_{\text{input}}(x) dx \end{bmatrix}^\top$$

where $\nabla A(\boldsymbol{\eta}_{1:2}) \equiv DA(\boldsymbol{\eta}_{1:2})^\top$ is the gradient vector of $A(\boldsymbol{\eta}_{1:2})$. It is then easy to show

$$\boldsymbol{\eta}_0^* = \log(C_{f_{\text{input}}}) - A(\boldsymbol{\eta}_{1:2}^*) - \frac{d}{2} \log(2\pi),$$

where

$$\boldsymbol{\eta}_{1:2}^* = (\nabla A)^{-1} \left\{ \begin{bmatrix} \int_{\mathbb{R}} x f_{\text{input}}(x)/C_{f_{\text{input}}} dx \\ \int_{\mathbb{R}} x^2 f_{\text{input}}(x)/C_{f_{\text{input}}} dx \end{bmatrix} \right\} \quad (2.40)$$

with existence and uniqueness of $(\nabla A)^{-1}$ being guaranteed by Proposition 3.2 of Wainwright & Jordan (2008),^[66] and $C_{f_{\text{input}}} \equiv \int_{\mathbb{R}} f_{\text{input}}(x) dx$. The Hessian matrix of $\mathcal{K}(\boldsymbol{\eta})$ is

$$\begin{aligned} \mathbf{HK}(\boldsymbol{\eta}) &= (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A(\boldsymbol{\eta}_{1:2})) \\ &\quad \times \left\{ \begin{bmatrix} 1 \\ \nabla A(\boldsymbol{\eta}_{1:2}) \end{bmatrix} \begin{bmatrix} 1 \\ \nabla A(\boldsymbol{\eta}_{1:2}) \end{bmatrix}^\top + \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{HA}(\boldsymbol{\eta}_{1:2}) \end{bmatrix} \right\}. \end{aligned}$$

By Proposition 3.1 of Wainwright & Jordan (2008),^[66] A is strictly convex on its domain and thus $\mathbf{HA}(\boldsymbol{\eta}_{1:2})$ is positive definite. Thus, $\mathbf{HK}(\boldsymbol{\eta})$ is positive definite for all $\boldsymbol{\eta}$ and so equation (2.40) is the unique maximiser of $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$. Therefore

$$\text{proj}_{\text{UN}}[f_{\text{input}}](x) = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta}^* \right\},$$

where $\boldsymbol{\eta}^*$ is as previously defined. However, $\boldsymbol{\eta}_{1:2}^*$ is the same natural parameter vector that arises via projection of $f_{\text{input}}/C_{f_{\text{input}}}$ onto the family of univariate normal density functions, thus

$$\text{proj}_{\text{N}}[f_{\text{input}}/C_{f_{\text{input}}}] (x) = \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix} \boldsymbol{\eta}_{1:2}^* - A(\boldsymbol{\eta}_{1:2}^*) \right\} (2\pi)^{-d/2}$$

which immediately leads to Result 3.

2.7.2 Proof of Result 5

We now detail closed form solutions to the integrals arising in the log-likelihood calculations of univariate GLMMs and its derivatives.

Lemma 1. For the following integrals of the form $\mathcal{W}_s(a, b) = \int_{-\infty}^{\infty} x^s \Phi(a + bx) \phi(x) dx$, where $s \in \{0, 1, 2, 3, 4, 5, 6\}$ and $a, b \in \mathbb{R}$, the following closed form solutions exist:

$$\mathcal{W}_0(a, b) = \int_{-\infty}^{\infty} \Phi(a + bx) \phi(x) dx = \Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \quad (2.41)$$

$$\mathcal{W}_1(a, b) = \int_{-\infty}^{\infty} x \Phi(a + bx) \phi(x) dx = \frac{b}{\sqrt{b^2 + 1}} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \quad (2.42)$$

$$\begin{aligned} \mathcal{W}_2(a, b) &= \int_{-\infty}^{\infty} x^2 \Phi(a + bx) \phi(x) dx = \Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) \\ &\quad - \frac{ab^2}{\sqrt{(b^2 + 1)^3}} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \end{aligned} \quad (2.43)$$

$$\begin{aligned} \mathcal{W}_3(a, b) &= \int_{-\infty}^{\infty} x^3 \Phi(a + bx) \phi(x) dx = \frac{|b| \left(2b^4 + (a^2 + 5)b^2 + 3\right)}{\sqrt{(b^2 + 1)^5}} \\ &\quad \times \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \end{aligned} \quad (2.44)$$

$$\begin{aligned} \mathcal{W}_4(a, b) &= \int_{-\infty}^{\infty} x^4 \Phi(a + bx) \phi(x) dx = 3\Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) \\ &\quad - \frac{ab^2 \left(3b^4 + (a^2 + 9)b^2 + 6\right)}{\sqrt{(b^2 + 1)^7}} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \end{aligned} \quad (2.45)$$

$$\begin{aligned} \mathcal{W}_5(a, b) &= \int_{-\infty}^{\infty} x^5 \Phi(a + bx) \phi(x) dx = (b^2 + 1)^{-9/2} \left(15 + 10(5 + a^2)b^2 \right. \\ &\quad \left. + (63 + 14a^2 + a^4)b^4 + 4(9 + a^2)b^6 + 8b^8\right) \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right), \end{aligned} \quad (2.46)$$

$$\begin{aligned} \mathcal{W}_6(a, b) &= \int_{-\infty}^{\infty} x^6 \Phi(a + bx) \phi(x) dx = 15\Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) - (b^2 + 1)^{-11/2} \\ &\quad \times \left\{ ab^2 \left(45 + 15(9 + a^2)b^2 + (150 + 20a^2 + a^4)b^4 + 5(15 + a^2)b^6 \right. \right. \\ &\quad \left. \left. + 15b^8 \right) \right\} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right). \end{aligned} \quad (2.47)$$

2.7.2.1 Proof of Lemma 1

We next provide the derivation of $\mathcal{W}_s(a, b)$ where $s = 3$. $s \in \{0, 1, 2, 4, 5, 6\}$ are solved analogously and hence their working out is omitted. For the integral

$$\mathcal{W}_s(a, b) = \int_{-\infty}^{\infty} x^3 \Phi(a + bx) \phi(x) dx,$$

consider the following conjecture:

$$\phi^{(3)}(x) \equiv (3x - x^3)\phi(x). \quad (2.48)$$

Proof: First note the derivatives

$$\begin{aligned}\phi'(x) &= -x\phi(x), \\ \phi''(x) &= (x^2 - 1)\phi(x), \\ \phi'''(x) &= (3x - x^3)\phi(x).\end{aligned}$$

Using simple algebraic manipulations

$$\begin{aligned}x^3\phi(x) &= -(3x - x^3)\phi(x) + 3x\phi(x) \\ &= -\phi'''(x) - 3\phi'(x).\end{aligned}$$

Thus

$$\begin{aligned}\int_{-\infty}^{\infty} x^3\Phi(a + bx)\phi(x)dx &= \int_{-\infty}^{\infty} \Phi(a + bx)\phi'''(x)dx - 3 \int_{-\infty}^{\infty} \Phi(a + bx)\phi'(x)dx \\ &= -T_A - 3T_B,\end{aligned}\tag{2.49}$$

where

$$T_A \equiv \int_{\mathbb{R}} \phi'''(x)\Phi(a + bx)dx \quad \text{and} \quad T_B \equiv \int_{\mathbb{R}} \phi'(x)\Phi(a + bx)dx.$$

Now

$$\begin{aligned}T_A &= \int_{\mathbb{R}} \phi'''(x)\Phi(a + bx)dx \\ &= -b \int_{\mathbb{R}} \phi''(x)\phi(a + bx)dx \\ &= - \int_{-\infty}^{\infty} \phi^{(2)}(x) \frac{1}{b} \phi^{(0)}\left(\frac{x - (-a/b)}{(1/b)}\right) dx \\ &= - \int_{-\infty}^{\infty} \phi^{(2)}(x) \phi_{\frac{1}{b}}^{(0)}(x - (-a/b)) dx.\end{aligned}$$

Then by C.1.12 of Wand & Jones (1995)^[67] it can be written

$$\begin{aligned}&= -(-1)^2 \phi_{\sqrt{1+\frac{1}{b^2}}}^{(2)}(-(-a/b)) \\ &= -\left(\frac{b^2}{1+b^2}\right)^{3/2} \phi^{(2)}\left(\frac{a}{\sqrt{1+b^2}}\right).\end{aligned}$$

Using the form of $\phi^{(2)}$ where x is given in the proof of the conjecture

$$\begin{aligned}&= -\frac{|b|^3}{\sqrt{(1+b^2)^3}} \left(\frac{a^2}{1+b^2} - 1\right) \phi\left(\frac{a}{\sqrt{(1+b^2)}}\right) \\ &= -\frac{|b|^3}{\sqrt{(1+b^2)^3}} \left(\frac{a^2 - (1+b^2)}{\sqrt{(1+b^2)^2}} - 1\right) \phi\left(\frac{a}{\sqrt{(1+b^2)}}\right) \\ &= -\frac{|b|^3(a^2 - b^2 - 1)}{\sqrt{(1+b^2)^5}} \phi\left(\frac{a}{\sqrt{(1+b^2)}}\right).\end{aligned}$$

The T_B part is calculated analogously to the first

$$\begin{aligned}
T_B &= \int_{-\infty}^{\infty} \phi'(x) \Phi(a + bx) dx \\
&= -b \int_{-\infty}^{\infty} \phi^{(0)}(x) \phi^{(0)}(a + bx) dx \\
&= - \int_{-\infty}^{\infty} \phi^{(0)}(x) \phi_{\frac{1}{b}}^{(0)}(x - (-a/b)) dx \\
&= - \left(1 + \frac{1}{b^2}\right)^{-1/2} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) \\
&= - \frac{|b|}{\sqrt{1 + b^2}} \phi\left(\frac{a}{\sqrt{b^2 + 1}}\right).
\end{aligned}$$

Substituting T_A and T_B into equation (2.49) leads directly to the required result. Two interesting patterns worth noting occur:

- The cumulative density function Φ occurs only for even values of s .
- The algebraic expression of the integral (equation (2.49) for the $s = 3$ case) takes a negative value when s is odd and positive when s is even.

Each moment of the target density shown in equation (2.5) can be obtained as

$$\mathcal{M}_k \equiv \int_{-\infty}^{\infty} x^k \Phi(c_0 + c_1 x) (2\pi)^{-1/2} \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} dx \exp(A(\boldsymbol{\eta})) (2\pi)^{1/2}.$$

Taking the inverse of natural parameters as per equation (1.5) yields,

$$\mathcal{M}_k = \int_{-\infty}^{\infty} x^k \Phi(c_0 + c_1 x) (2\pi)^{-1/2} \exp \left\{ -\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right\} dx \exp(A(\boldsymbol{\eta})) (2\pi)^{1/2} \sigma.$$

Let $u = \frac{x - \mu}{\sigma}$, then $x = \mu + \sigma u$ and $dx = \sigma du$, then

$$\mathcal{M}_k = \int_{-\infty}^{\infty} (\mu + \sigma u)^k \Phi(c_0 + c_1 \mu + c_1 \sigma u) \phi(u) du Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + 1/2 \log(2\pi) + \log(\sigma))$. Using Lemma 1, it is easy to find each of the k th moments

$$\begin{aligned}
Z_1^{-1} \mathcal{M}_0 &= \Phi(r_2), \\
Z_1^{-1} \mathcal{M}_1 &= \mu \Phi(r_2) + (2c_1 \sigma^2 r_1^{-1}) \phi(r_2), \\
Z_1^{-1} \mathcal{M}_2 &= \mu^2 \Phi(r_2) + (4c_1 \sigma^2 \mu r_1^{-1}) \phi(r_2) + \sigma^2 \Phi(r_2) - (4r_2 c_1^2 \sigma^4 r_1^{-2}) \phi(r_2),
\end{aligned}$$

for $r_2 = 2(c_0 + c_1 \mu) r_1^{-1}$ and $r_1 = 2\sqrt{(c_1 \sigma)^2 + 1}$. The optimal mean and variance parameters follow respectively,

$$\mu^* = E(x) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \mu + (2c_1 \sigma^2 r_1^{-1}) \zeta'(r_2)$$

and

$$\begin{aligned} (\sigma^*)^2 = E(x^2) - E(x)^2 &= \frac{\mathcal{M}_2}{\mathcal{M}_0} - \left(\frac{\mathcal{M}_1}{\mathcal{M}_0}\right)^2 \\ &= \sigma^2 - (4c_1^2\sigma^4r_1^{-2})(r_2 - \zeta'(r_2))\zeta'(r_2). \end{aligned}$$

By converting to natural parameters and using matrix notation we arrive at the required result. Stable computation of $\zeta'(r_2)$ and $\zeta''(r_2)$ in R is provided by the function `zeta()` in the package “sn” (Azzalini, 2016^[4]), where the derivatives are controlled by the argument “k”.

2.7.3 Proof of Result 7

We next provide the derivation of $E(x^s)$ for $s = 2$. Note $s \in \{0, 1, 3, 4, 5, 6\}$ are solved analogously and hence their working out is omitted. First note, for each s the integral

$$\int_{-\infty}^{\infty} x^s \phi(x) dx$$

takes the following values:

$$\begin{aligned} \int_{-\infty}^{\infty} x^0 \phi(x) dx &= 1, & \int_{-\infty}^{\infty} x^1 \phi(x) dx &= 0, & \int_{-\infty}^{\infty} x^2 \phi(x) dx &= 1, \\ \int_{-\infty}^{\infty} x^3 \phi(x) dx &= 0, & \int_{-\infty}^{\infty} x^4 \phi(x) dx &= 3, & \int_{-\infty}^{\infty} x^5 \phi(x) dx &= 0, \\ \int_{-\infty}^{\infty} x^6 \phi(x) dx &= 15. \end{aligned}$$

With this in mind consider

$$\begin{aligned} E(u^2) &= \int_{-\infty}^{\infty} u^2 \exp(T(u)^\top \eta - A(\eta)) (2\pi)^{-1/2} du \\ &= \int_{-\infty}^{\infty} u^2 \frac{1}{\sigma} \phi\left(\frac{u - \mu}{\sigma}\right) du. \end{aligned}$$

Then employing the change of variables $z = \frac{u - \mu}{\sigma} \implies u = \mu + \sigma z$, it can be written

$$\begin{aligned} E(u^2) &= \int_{-\infty}^{\infty} (\mu + \sigma z)^2 \phi(z) dz \\ &= \mu^2 \int_{-\infty}^{\infty} \phi(z) dz + 2\mu\sigma \int_{-\infty}^{\infty} z \phi(z) dz + \sigma^2 \int_{-\infty}^{\infty} z^2 \phi(z) dz. \end{aligned}$$

Noting $\int_{-\infty}^{\infty} z \phi(z) dz$ integrates to zero while $\int_{-\infty}^{\infty} \phi(z) dz$ and $\int_{-\infty}^{\infty} z^2 \phi(z) dz$ integrates to 1 leads to the required result.

2.7.4 Details on finding the inverse function gradient map log-partition function for each s

The following section of appendix aims to more explicitly explain the function $(\nabla A)^{-1}(\boldsymbol{\tau}, s)$ and the calculations it involves for each s . We only provide a full explanation for the $s = 0$ case since the other cases are analogous.

2.7.4.1 The $s = 0$ case

Firstly note that

$$E(x^0) = 1.$$

Following from the previous section, $E(x^0)$ can be used to normalise the approximation of the input density

$$f_{\text{UN}}(x) = (2\pi\sigma^2)^{-1/2} \exp\left(- (x - \mu)^2 / (2\sigma^2)\right).$$

Next note that using the standard change of variables and natural parameters

$$\begin{aligned} f_{\text{UN}}(x) &= (2\pi)^{-1/2} \exp\left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta}_{-1} - \left(-\eta_1^2 / (4\eta_2) - 1/2 \log(-2\eta_2) \right) \right\} \\ &= \exp\left(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1} - A(\boldsymbol{\eta}_{-1}) \right) h(x), \end{aligned}$$

where

$$h(x) \equiv (2\pi)^{-1/2}, \quad A(\boldsymbol{\eta}_{-1}) = -\eta_1^2 / (4\eta_2) - 1/2 \log(-2\eta_2).$$

The derivative vector of the log-partition function can be shown to be

$$(\nabla A)(\boldsymbol{\eta}_{-1}) = \begin{bmatrix} -\frac{\eta_1}{2\eta_2} \\ \left(\frac{\eta_1}{2\eta_2}\right)^2 - \frac{1}{2\eta_2} \end{bmatrix}.$$

We can now use the derivative vector of the log-partition to find an estimate of the minimum sum of squares deviation from the root of the inverse-defining equation using a Nelder-Mead search implemented through the `optim()` function in R. We then obtain the Hessian of the derivative vector of the log-partition function

$$\mathbf{H}(A(\boldsymbol{\eta}_{-1})) = \begin{bmatrix} -\frac{1}{2\eta_2} & \frac{\eta_1}{2\eta_2^2} \\ \frac{\eta_1}{2\eta_2^2} & \frac{\eta_2 - \eta_1^2}{2\eta_2^3} \end{bmatrix}$$

such that we can implement a Newton Raphson search to further refine the approximation of $\boldsymbol{\eta}^*$ using the minimum calculated through our Nelder-Mead search as a starting value.

2.7.4.2 The $s = 2$ case

Following calculations analogous to the $s = 0$ case it is easy to show:

$$(\nabla A)(\boldsymbol{\eta}_{-1}) = \begin{bmatrix} \frac{\eta_1(6\eta_2 - \eta_1^2)}{2\eta_2(\eta_1^2 - 2\eta_2)} \\ \frac{12\eta_1^2\eta_2 - 12\eta_2^2 - \eta_1^4}{4\eta_2^2(2\eta_2 - \eta_1^2)} \end{bmatrix}$$

and

$$H(A(\boldsymbol{\eta}_{-1})) = \begin{bmatrix} \frac{-\eta_1^4 - 12\eta_2^2}{2\eta_2(\eta_1^2 - 2\eta_2)^2} & \frac{\eta_1(12\eta_2^2 - 4\eta_1^2\eta_2 + \eta_1^4)}{2\eta_2^2(2\eta_2 - \eta_1^2)^2}, \\ \frac{\eta_1(12\eta_2^2 - 4\eta_1^2\eta_2 + \eta_1^4)}{2\eta_2^2(2\eta_2 - \eta_1^2)^2} & \frac{12\eta_2^3 - 24\eta_1^2\eta_2^2 + 9\eta_1^4\eta_2 - \eta_1^6}{2\eta_2^3(2\eta_2 - \eta_1^2)^2} \end{bmatrix}.$$

2.7.4.3 The $s = 4$ case

As in the first derivative case, the calculations required are analogous to the $s = 0$ case, and so

$$(\nabla A)(\boldsymbol{\eta}_{-1}) = \begin{bmatrix} \frac{4\eta_1^3 - 24\eta_2\eta_1}{12\eta_2^2 - 12\eta_2\eta_1^2 + \eta_1^4} - \frac{\eta_1}{2\eta_2} \\ \frac{180\eta_1^2\eta_2^2 - 120\eta_2^3 - 30\eta_1^4\eta_2 + \eta_1^6}{4\eta_2^2(12\eta_2^2 - 12\eta_1^2\eta_2 + \eta_1^4)} \end{bmatrix}$$

and

$$H(A(\boldsymbol{\eta}_{-1})) = \begin{bmatrix} \frac{16\eta_2\eta_1^6 - \eta_1^8 - 120\eta_2^2\eta_1^4 - 720\eta_2^4}{2\eta_2(12\eta_2^2 - 12\eta_1^2\eta_2 + \eta_1^4)^2} & \frac{\eta_1(720\eta_2^4 - 480\eta_1^2\eta_2^3 + 216\eta_1^4\eta_2^2 - 24\eta_1^6\eta_2 + \eta_1^8)}{2\eta_2^2(12\eta_2^2 - 12\eta_1^2\eta_2 + \eta_1^4)^2} \\ \frac{\eta_1(720\eta_2^4 - 480\eta_1^2\eta_2^3 + 216\eta_1^4\eta_2^2 - 24\eta_1^6\eta_2 + \eta_1^8)}{2\eta_2^2(12\eta_2^2 - 12\eta_1^2\eta_2 + \eta_1^4)^2} & \frac{720\eta_2^5 - 2160\eta_1^2\eta_2^4 + 1560\eta_1^4\eta_2^3 - 384\eta_1^6\eta_2^2 + 33\eta_1^8\eta_2 - \eta_1^{10}}{2\eta_2^3(12\eta_2^2 - 12\eta_1^2\eta_2 + \eta_1^4)^2} \end{bmatrix}$$

Chapter 3

Expectation propagation for general one level probit mixed models

With knowledge from Chapter 2 on the random intercepts only model, we now extend our model to the more general case of GLMMs discussed in Section 1.7 that allow for any number of fixed and random effects. We aim to find an approximation to the maximum likelihood of the parameters $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ with 95% confidence intervals. Our model is specified as in equation (3.1)

$$y_{ij} | \mathbf{u}_i \stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left(\Phi(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right), \quad \mathbf{u}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}_{d_{\mathbf{R}}}, \boldsymbol{\Sigma}),$$
$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad (3.1)$$

where the notation follows from the general one level model presented in Section 1.8. The log-likelihood can be expressed as

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}),$$

where

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \log \int_{\mathbb{R}^{d_{\mathbf{R}}}} \left\{ \prod_{j=1}^{n_i} \Phi \left((2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right) \right\}$$
$$\times |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right) d\mathbf{u}_i \quad (3.2)$$

and the best predictor of \mathbf{u}_i is

$$\text{BP}(\mathbf{u}_i) \equiv \frac{\int_{\mathbb{R}^{d^{\mathbf{R}}}} \mathbf{u} \left\{ \prod_{j=1}^{n_i} \Phi \left((2y_{ij} - 1) (\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right) \right\} \exp \left(-\frac{1}{2} \mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) d\mathbf{u}}{\int_{\mathbb{R}^{d^{\mathbf{R}}}} \left\{ \prod_{j=1}^{n_i} \Phi \left((2y_{ij} - 1) (\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right) \right\} \exp \left(-\frac{1}{2} \mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{u} \right) d\mathbf{u}}. \quad (3.3)$$

We are now left with the more complex problem of intractable $d^{\mathbf{R}}$ -dimensional integrals. As mentioned in Section 2.1, although Gauss-Hermite quadrature can be used to obtain high accuracy results for models with one random effect (i.e. $d^{\mathbf{R}} = 1$), it is often not feasible computationally for the case where more than two or three random effects exist. It is in this case that the benefits of approximate inference, specifically via EP, become more obvious. Additionally, the maximisation problem of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ is complicated over a changing dimensional space, so we utilise the quasi-Newton approach discussed in the Chapter 2.

The structure of this chapter follows that of the previous, first explaining the schematic of likelihood approximation using EP in Section 3.1, before discussing computation of point estimates and confidence intervals in Section 3.2. The calculation of best predictors for the random effects is given in Section 3.3 before the results of our simulation studies are presented in Section 3.4.

3.1 Expectation propagation likelihood approximation

We now show how EP can be used to approximate the likelihood by updating and summing the natural parameter updates. Details of the required projections and the message passing formulation to compartmentalise the algebra are provided. As in the univariate case, we can approximate each $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ and sum them to obtain the full log-likelihood. EP is motivated by the selection of an unnormalised multivariate normal density function that minimises the KL divergence criterion to replace each

$$\Phi \left((2y_{ij} - 1) (\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right), \quad 1 \leq j \leq n_i$$

in equation (3.2). Doing so for the probit case constrains the integrand to be a product of multivariate normal density functions with closed form solutions.

Consider the full KL divergence presented in equation (1.19) and the following family of unnormalised multivariate normal density functions written in exponential family

form

$$f_{\text{UN}}(\mathbf{x}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \begin{bmatrix} \eta_0 \\ \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix} \right\}, \quad (3.4)$$

with natural parameters $\eta_0 \in \mathbb{R}$, $\boldsymbol{\eta}_1$ (a $d \times 1$ vector) and $\boldsymbol{\eta}_2$ (a $\frac{1}{2}d(d+1)$ vector). In the higher dimensional general case of GLMM, the goal of the EP problem is to find the optimal natural parameters (denoted by η_0^* , $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$) which minimise $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$, where $f_{\text{input}} \in L_1(\mathbb{R}^d)$. This solution is referred to as the KL projection onto the family of multivariate normal density functions and is written as

$$\text{proj}_{\mathcal{N}}[f_{\text{input}}](\mathbf{x}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta}^* \right\}, \quad (3.5)$$

where

$$\boldsymbol{\eta}^* \equiv \begin{bmatrix} \eta_0^* \\ \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}$$

and

$$(\eta_0^*, \boldsymbol{\eta}_1^*, \boldsymbol{\eta}_2^*) = \underset{(\eta_0, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2) \in H}{\text{argmin}} \text{KL}(f_{\text{input}} \parallel f_{\text{UN}}).$$

In the special case of KL projection onto the unnormalised multivariate normal family, this problem simplifies further to moment-matching, where $(\eta_0^*, \boldsymbol{\eta}_1^*, \boldsymbol{\eta}_2^*)$ is the unique vector that matches the zeroth, first and second order moments of f_{UN} and f_{input} .

For probit binary GLMMs, EP requires repeated projection of

$$f_{\text{input}}(\mathbf{x}) = \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp \left((\boldsymbol{\eta}_1^{\text{input}})^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x} \right) \quad (3.6)$$

onto an unnormalised multivariate normal distribution, $c_0 = (2y_{ij} - 1)\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}$, $\mathbf{c}_1 =$

$(2y_{ij} - 1)\mathbf{x}_{ij}^{\mathbf{R}}$ and $\mathbf{x} = \mathbf{u}_i$. With this in mind, we seek $\boldsymbol{\eta}^*$ such that

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp \left\{ \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta}^{\text{input}} \right\} d\mathbf{x} \\ = \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta}^* \right\} d\mathbf{x}, \end{aligned} \quad (3.7)$$

where $k \in \{0, 1, 2\}$.

Now consider the multivariate extension of Result [3](#).

Result 12. For an unnormalised input function $f \in L_1(\mathbb{R}^d)$ such that $f \geq 0$ for all $x \in \mathbb{R}$ where $C_f \equiv \int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x}$, the projection onto the unnormalised multivariate normal family is

$$\text{proj}_{UN}[f](\mathbf{x}) = C_f \text{proj}_{\mathbf{N}}[f/C_f](\mathbf{x}),$$

where $\text{proj}_{\mathbf{N}}$ is the projection onto the multivariate normal family.

Obtaining the natural parameters $\boldsymbol{\eta}^*$ for projection onto the unnormalised multivariate normal family follows from obtaining the projection onto the multivariate normal family. More explicitly, the optimal natural parameters $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$, are given according to the projection of the normalised function $f_{\text{input}}/C_{f_{\text{input}}}$ onto the multivariate normal family. We can subsequently use these optimal natural parameters to find the normalising natural parameter $\boldsymbol{\eta}_0^*$ via Result [13](#) and thus obtain the projection onto unnormalised normal family.

Result 13. When f_{input} density follows the form of equation [\(2.5\)](#), $\boldsymbol{\eta}_0^*$ is given by

$$\boldsymbol{\eta}_0^* = \log(C_{f_{\text{input}}}) - A(\boldsymbol{\eta}_1^*, \boldsymbol{\eta}_2^*) - \frac{d}{2} \log(2\pi),$$

where the log-partition function is as defined in Section [1.5.2.2](#).

Thus to obtain the required projection, we first obtain the optimal natural parameters $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$ to project onto the multivariate normal family as is presented in Result [14](#).

Result 14. Given f_{input} follows the form of equation (3.6), the projection onto the multivariate normal family is given by

$$\text{proj}_{\mathcal{N}}[f_{input}] = \exp(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*))h(\mathbf{x}),$$

where

$$\boldsymbol{\eta}_{-1}^* = K_{\text{probit}}(\boldsymbol{\eta}_{-1}^{input}; c_0, \mathbf{c}_1), \quad \boldsymbol{\eta}_{-1}^{input} \equiv \begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix},$$

$K_{\text{probit}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right)$ is defined in Definition 13 and $\mathbf{T}(\mathbf{x})$ and $h(\mathbf{x})$ follow from Section 1.5.2.2

Definition 13. For primary arguments \mathbf{a}_1 ($d \times 1$) and \mathbf{a}_2 ($\frac{1}{2}d(d+1) \times 1$) such that $\text{vec}^{-1}(-\mathbf{D}_d^{+\top} \mathbf{a}_2)$ is symmetric and positive definite, and auxiliary arguments $c_0 \in \mathbb{R}$ and \mathbf{c}_1 ($d \times 1$), the function $K_{\text{probit}} : H \rightarrow H$ is given by

$$K_{\text{probit}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \equiv \begin{bmatrix} \mathbf{R}_5^\top (\mathbf{a}_1 + r_3 \mathbf{c}_1) \\ \mathbf{D}_d^\top \text{vec}(\mathbf{R}_5^\top \mathbf{A}_2) \end{bmatrix}, \quad (3.8)$$

where

$$\begin{aligned} \mathbf{A}_2 &\equiv \text{vec}^{-1}(\mathbf{D}_d^{+\top} \mathbf{a}_2), \quad r_1 \equiv \sqrt{2(2 - \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{c}_1)}, \quad r_2 \equiv (2c_0 - \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1) r_1^{-1} \\ r_3 &\equiv 2\zeta'(r_2) r_1^{-1}, \quad r_4 \equiv -2\zeta''(r_2) r_1^{-2} \quad \text{and} \quad \mathbf{R}_5 \equiv (\mathbf{A}_2 + r_4 \mathbf{c}_1 \mathbf{c}_1^\top)^{-1} \mathbf{A}_2. \end{aligned}$$

The proof of Definition 13 is given in Appendix 3.5.1. Using Result 14 we now obtain the normalising natural parameter $\boldsymbol{\eta}_0^*$ to find the projection onto unnormalised normal family.

3.1.1 Projection onto the unnormalised multivariate normal family

Recall the moment matching problem described by equation (3.7) and Result 13. Then, as in the univariate case we require

$$C_{f_{\text{input}}} = \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} = (2\pi)^{d/2} \exp(A(\boldsymbol{\eta}^{\text{input}})) \Phi(r_2),$$

where r_2 is given in Definition 13 and $A(\boldsymbol{\eta})$ is defined in Section 1.5.2.2. Analogous to the argument given in the univariate case we then get

$$\boldsymbol{\eta}_0^* = \log \Phi(r_2) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^\top (\mathbf{H}_2^*)^{-1} \boldsymbol{\eta}_1^* - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^\top (\mathbf{H}_2^{\text{input}})^{-1} \boldsymbol{\eta}_1^{\text{input}} + \frac{1}{2} \log (|\mathbf{H}_2^*|/|\mathbf{H}_2^{\text{input}}|).$$

Definition 14. Consider first, primary arguments \mathbf{a}_1 and \mathbf{b}_1 and auxiliary argument \mathbf{c}_1 where all three are $d \times 1$. Next consider arguments \mathbf{a}_2 and \mathbf{b}_2 which are all $(\frac{1}{2}d(d+1) \times 1)$ such that both $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{a}_2)$ and $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{b}_2)$ are symmetric and positive definite. Finally note auxiliary argument $c_0 \in \mathbb{R}$. Then the function $C_{\text{probit}} : H \times H \rightarrow \mathbb{R}$ is given by

$$C_{\text{probit}} \left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}; c_0, \mathbf{c}_1 \right) \equiv \log \Phi(r_2) + \frac{1}{4} \mathbf{b}_1^\top \mathbf{B}_2^{-1} \mathbf{b}_1 - \frac{1}{4} \mathbf{a}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1 + \frac{1}{2} \log (|\mathbf{B}_2|/|\mathbf{A}_2|),$$

where $\mathbf{A}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{a}_2)$, $\mathbf{B}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{b}_2)$ and r_2 follows from Definition 13.

In summary, the calculations required obtain the KL projection of the input function onto the unnormalised multivariate normal family are given by Result 15.

Result 15. For an unnormalised input function of the form of equation (3.6) then

$$\text{proj}_{UN}[f_{input}](x) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} \right\},$$

where

$$\begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} = K_{\text{probit}} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix}; c_0, \mathbf{c}_1 \right) \quad \text{and} \quad \eta_0^* = C_{\text{probit}} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1 \right).$$

We now show how this result can be used in a message passing framework to minimise the algebra required.

3.1.2 Message passing formulation

Note the components of the sum of the log-likelihood function $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ can be written as

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \log \int_{\mathbb{R}^{d\mathbf{R}}} \left(\prod_{j=1}^{n_i} p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) d\mathbf{u}_i, \quad (3.9)$$

where

$$p(y_{ij} | \mathbf{u}_i, \boldsymbol{\beta}) \equiv \Phi \left((2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right)$$

and

$$p(\mathbf{u}_i; \boldsymbol{\Sigma}) \equiv |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right)$$

are respectively the conditional density functions of each response given its random effect and the density function of that random effect. Also note the alternate matrix expression

$$p(\mathbf{u}_i; \boldsymbol{\Sigma}) \equiv \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}_\Sigma \right\}, \quad \text{where} \quad \boldsymbol{\eta}_\Sigma \equiv \begin{bmatrix} -\frac{1}{2} \log |2\pi\boldsymbol{\Sigma}| \\ 0_{d\mathbf{R}} \\ -\frac{1}{2} \mathbf{D}_{d\mathbf{R}}^\top \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}. \quad (3.10)$$

The structure of the product of equation (3.9) can be represented as Figure 3.1, where circular stochastic nodes correspond to the random vector \mathbf{u}_i , solid squares indicate the $n_i + 1$ factor nodes and the dependencies of the factor nodes on the stochastic node \mathbf{u}_i is demonstrated through the edges.

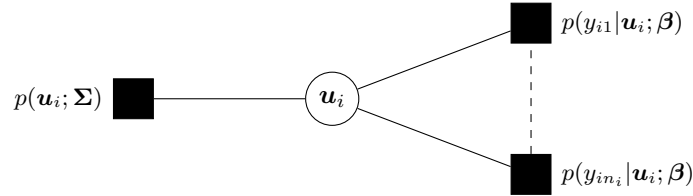


Figure 3.1: Factor graph representation of the product structure of the integrand in equation (3.9).

We now proceed in a manner analogous to Section 2.2.2 for the univariate case, using the Bayesian approach of Minka (2005).⁴⁴ The EP approximation of $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ involves projection onto the unnormalised multivariate normal family. Suppose that

$$p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) = \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{array} \right]^\top \boldsymbol{\eta}_{ij} \right\}, \quad 1 \leq j \leq n_i,$$

is initialised to be from the family of unnormalised multivariate normal density functions in \mathbf{u}_i . Then, for each $j = 1, \dots, n_i$, the $\boldsymbol{\eta}_{ij}$ update involves minimisation of

$$\text{KL} \left(p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \left(\prod_{j' \neq j}^{n_i} p(y_{ij'} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) \parallel \left(\prod_{j'=1}^{n_i} p(y_{ij'} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) \right) \quad (3.11)$$

as a function of \mathbf{u}_i . This can be achieved using Result 15 to update $\boldsymbol{\eta}_{ij}$ in an iterative procedure until it converges.

Using message passing, we compartmentalise the otherwise cumbersome algebra via a simple extension of the steps involved in the univariate case. As such we do not repeat all of the steps again, instead providing the main key equations, which mirror equation (54) and (83) of Minka (2005).⁴⁴ Since derivation of the simplification required for the structure of the message passing scheme is a simple extension of the univariate case, we do not repeat it. The KL divergence of equation (3.11) can be re-expressed, such that

the messages from the factor $p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ to stochastic node \mathbf{u}_i are updated according to

$$m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}_{\text{UN}} [m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})](\mathbf{u}_i)}{m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i)}, \quad 1 \leq j \leq n_i \quad (3.12)$$

and the update of message passed from $p(\mathbf{u}_i; \boldsymbol{\Sigma})$ to \mathbf{u}_i is

$$m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}_{\text{UN}} [m_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \boldsymbol{\Sigma})}(\mathbf{u}_i) p(\mathbf{u}_i; \boldsymbol{\Sigma})](\mathbf{u}_i)}{m_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \boldsymbol{\Sigma})}(\mathbf{u}_i)}. \quad (3.13)$$

Similarly, the updates of stochastic node to factor messages are

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})}(\mathbf{u}_i) = m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \prod_{j' \neq j}^{n_i} m_{p(y_{ij'}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i), \quad 1 \leq j \leq n_i \quad (3.14)$$

and

$$m_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \boldsymbol{\Sigma})}(\mathbf{u}_i) = \prod_{j=1}^{n_i} m_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i). \quad (3.15)$$

We now seek any algebraic simplifications of the key messages, particularly use of natural parameters. Recall that $p(\mathbf{u}_i; \boldsymbol{\Sigma})$ can be written using natural parameters in the form of equation (3.10) and that the unnormalised normal density constraint is enforced on equations (3.12) and (3.14). Then

$$m_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \boldsymbol{\Sigma})}(\mathbf{u}_i) = \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{array} \right]^\top \boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(\mathbf{u}_i; \boldsymbol{\Sigma})} \right\}. \quad (3.16)$$

Substituting the above forms into equation (3.13) leads to

$$m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow p(\mathbf{u}_i; \boldsymbol{\Sigma}) = \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{array} \right]^\top \boldsymbol{\eta}_{\boldsymbol{\Sigma}} \right\}.$$

This implies the message $m_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ is constant throughout the message passing updates. As such, we now set

$$\boldsymbol{\eta}_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i} \leftarrow \boldsymbol{\eta}_{\boldsymbol{\Sigma}}. \quad (3.17)$$

For convenience, we denote the natural parameter vector

$$\boldsymbol{\eta}^{\otimes} \equiv \boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}.$$

Recall that from equation (3.16)

$$\begin{aligned} m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}(\mathbf{u}_i) &= \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}^{\otimes} \right\} \\ &= \exp(\eta_0^{\otimes}) \exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^{\otimes} + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^{\otimes} \right\}. \end{aligned}$$

Substituting this into equation (3.12) and following simplifications analogous to the univariate case lead to

$$m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \longleftarrow \frac{\text{proj}_{\mathcal{UN}} \left[\Phi(c_0 + \mathbf{c}_{1ij}^\top \mathbf{u}_i) \exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^{\otimes} + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^{\otimes} \right\} \right]}{\exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^{\otimes} + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^{\otimes} \right\}},$$

where $c_0 = (2y_{ij} - 1)\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{c}_{1ij} \equiv (2y_{ij} - 1)\mathbf{x}_{ij}^{\mathbf{R}}$. Utilising Result 15,

$$m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \longleftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right\}, \quad (3.18)$$

where the linear and quadratic coefficient updates are

$$(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{1:2} \longleftarrow K_{\text{probit}}(\boldsymbol{\eta}_{1:2}^{\otimes}; c_0, \mathbf{c}_{1ij}) - \boldsymbol{\eta}_{1:2}^{\otimes} \quad (3.19)$$

and the constant coefficient update is

$$(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_0 \longleftarrow C_{\text{probit}}(\boldsymbol{\eta}_{1:2}^{\otimes}, (\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{1:2} + \boldsymbol{\eta}_{1:2}^{\otimes}; c_0, \mathbf{c}_{1ij}).$$

Using the simplification of equation (3.12) and (3.13), equation (3.14) can be shown to be

$$m_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}(\mathbf{u}_i) \longleftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})} \right\}, \quad (3.20)$$

where

$$\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})} \leftarrow \boldsymbol{\eta}_{p(\mathbf{u}_i;\boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i} + \sum_{j' \neq j} \boldsymbol{\eta}_{p(y_{ij'}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i}.$$

The EP approximation of each component $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ of the log-likelihood sum is given by

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \log \int_{\mathbb{R}^{d\mathbf{R}}} \left(\prod_{j=1}^{n_i} m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \right) m_{p(\mathbf{u}_i;\boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) d\mathbf{u}_i. \quad (3.21)$$

The success of EP depends on each of the messages in equation (3.21) being an unnormalised multivariate normal density function and that a closed form solution to the integral as follows:

$$\begin{aligned} & \int_{\mathbb{R}^{d\mathbf{R}}} \left(\prod_{j=1}^{n_i} m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \right) m_{p(\mathbf{u}_i;\boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) d\mathbf{u}_i \\ &= \int_{\mathbb{R}^{d\mathbf{R}}} \prod_{j=1}^{n_i} \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right\} \\ & \quad \times \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\mathbf{u}_i;\boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i} \right\} d\mathbf{u}_i \\ &= (2\pi)^{d\mathbf{R}/2} \exp \left\{ \left(\boldsymbol{\eta}_{\boldsymbol{\Sigma}} + \sum_{j=1}^{n_i} \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_0 + A \left(\left(\boldsymbol{\eta}_{\boldsymbol{\Sigma}} + \sum_{j=1}^{n_i} \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{1:2} \right) \right\}. \end{aligned}$$

The full algorithm for the approximation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ using EP is provided in Algorithm 6.

Algorithm 6 *Explicit form of algorithm used for the message passing approach to EP*

Inputs: y_{ij} , $\mathbf{x}_{ij}^{\mathbf{F}}$, $\mathbf{x}_{ij}^{\mathbf{R}}$, $1 \leq i \leq m$, $1 \leq j \leq n_i$; $\boldsymbol{\beta}$ ($d^{\mathbf{F}} \times 1$), $\boldsymbol{\Sigma}$ ($d^{\mathbf{R}} \times d^{\mathbf{R}}$, are symmetric and positive definite).

Set constants: $c_{0_{ij}} \leftarrow (2y_{ij} - 1)(\boldsymbol{\beta}^{\top} \mathbf{x}_{ij}^{\mathbf{F}})$, $c_{1_{ij}} \leftarrow (2y_{ij} - 1)\mathbf{x}_{ij}^{\mathbf{R}}$, $1 \leq i \leq m$, $1 \leq j \leq n_i$.

$$\boldsymbol{\eta}_{p(\mathbf{u}_i; \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i} \leftarrow \boldsymbol{\eta}_{\boldsymbol{\Sigma}} \equiv \begin{bmatrix} -\frac{1}{2} \log |2\pi \boldsymbol{\Sigma}| \\ \mathbf{0}_{d^{\mathbf{R}}} \\ -\frac{1}{2} \mathbf{D}_{d^{\mathbf{R}}}^{\top} \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}, \quad 1 \leq i \leq m.$$

For $i = 1, \dots, m$:

Initialise: $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$, $1 \leq j \leq n_i$ as per the equation (3.24).

Cycle:

$$\text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) \leftarrow \sum_{j=1}^{n_i} \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$$

For $j = 1, \dots, n_i$:

$$\begin{aligned} \boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{u}_i, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_i} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) - \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \\ (\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{1:2} &\leftarrow K_{\text{probit}} \left((\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{1:2}; c_{0_{ij}}, c_{1_{ij}} \right) \\ &\quad - (\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{1:2} \end{aligned}$$

until convergence of all natural parameters vectors.

For $j = 1, \dots, n_i$:

$$\begin{aligned} (\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_0 &\leftarrow C_{\text{probit}} \left((\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{1:2}, \right. \\ &\quad \left. (\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})_{1:2} + (\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})})_{1:2}; c_{0_{ij}}, c_{1_{ij}} \right) \end{aligned}$$

$$\text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) \leftarrow \sum_{j=1}^{n_i} \boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}.$$

Output: The full approximate log-likelihood is given by

$$\begin{aligned} \underline{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \frac{md^{\mathbf{R}}}{2} \log(2\pi) + \sum_{i=1}^m \left\{ \left(\boldsymbol{\eta}_{\boldsymbol{\Sigma}} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) \right)_0 \right. \\ &\quad \left. + A \left\{ \left(\boldsymbol{\eta}_{\boldsymbol{\Sigma}} + \text{SUM}(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) \right)_{1:2} \right\} \right\} \end{aligned}$$

where, $A(\boldsymbol{\eta})$ is defined as in equation (1.7) and $\boldsymbol{\eta}_{\boldsymbol{\Sigma}}$ follows from equation (3.10).

3.1.3 Starting values for Algorithm 6

The EP message passing algorithm proposed relies on good starting values for convergence. We now derive starting values for $\eta_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ using a Taylor series expansion. Note that

$$\log p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) = \zeta(a_{ij}) - \log(2), \quad \text{where} \quad a_{ij} \equiv (2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})$$

and ζ is defined as before. Let $\hat{\mathbf{u}}_i$ be a Laplace approximation to \mathbf{u}_i . Now consider the following Taylor series expansion of the data dependent component of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$:

$$\begin{aligned} \zeta(a_{ij}) &= \zeta(\hat{a}_{ij} + (2y_{ij} - 1)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}}) \\ &= \zeta(\hat{a}_{ij}) + (2y_{ij} - 1)\zeta'(\hat{a}_{ij})(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}} + \frac{1}{2}((\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}})^2 \zeta''(\hat{a}_{ij}) + \dots \\ &= \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top) \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where $\hat{a}_{ij} \equiv (2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \hat{\mathbf{u}}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})$ and

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} \zeta(\hat{a}_{ij}) \\ (2y_{ij} - 1)\zeta'(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}} \\ \frac{1}{2}\zeta''(\hat{a}_{ij})\mathbf{D}_{d\mathbf{R}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix}.$$

It follows that the quadratic approximation to $\log p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ based on Taylor expansion about $\hat{\mathbf{u}}_i$ is $\log \check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ where

$$\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) \equiv \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top) \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} \right\}. \quad (3.22)$$

The starting value recommendation for $\eta_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ is based on replacement of $p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ by $\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ in equation (3.17):

$$m_{\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}[m_{\mathbf{u}_i\rightarrow\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}(\mathbf{u}_i)\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})](\mathbf{u}_i)}{m_{\mathbf{u}_i\rightarrow\check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}(\mathbf{u}_i)} = \check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}). \quad (3.23)$$

Note that in this case, since $\check{p}(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta})$ is already univariate normal the projection is superfluous. The starting values for $\eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$ that arises from this substitution is

$$\exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{bmatrix}^\top \eta_{p(y_{ij}|\mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}^{\text{start}} \right\} = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top) \end{bmatrix}^\top \check{\eta}_{ij} \right\}.$$

By matching coefficients of like terms it is easy to show

$$\eta_{p(y_{ij}|\mathbf{u}_i) \rightarrow \mathbf{u}_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ (2y_{ij} - 1)\zeta'(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}} - \zeta''(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i \\ \frac{1}{2}\zeta''(\hat{a}_{ij})\mathbf{D}_{d^{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix}, \quad (3.24)$$

where

$$\eta_0^{\text{start}} = \zeta(\hat{a}_{ij}) - (2y_{ij} - 1)\zeta'(\hat{a}_{ij})(\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i + \frac{1}{2}\zeta''(\hat{a}_{ij})((\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i)^2.$$

In Algorithm [3](#) η_0^{start} is not used in the cycle loop and thus can be set to any arbitrary number without affecting the algorithm. We use Laplace approximation to estimate $\hat{\mathbf{u}}_i$. For the R computing environment, the function `glmer()` of the package “lme4” (Bolker, et al., 2018[5](#)) provides fast Laplace approximation-based predictions for the \mathbf{u}_i .

3.2 Computation of point estimates and confidence intervals

We now address the computations required to find point estimates and confidence intervals for the parameters $\boldsymbol{\beta}, \boldsymbol{\Sigma}$. Although we still use the R function `optim()` in the “stats” package (R Core Team, 2019[56](#)) the higher dimensionality introduces new problems, particularly when optimising the variance parameter $\boldsymbol{\Sigma}$ since it is constrained to be symmetric and positive definite. Before conducting any optimisation, we must ensure the search occurs over the cone of symmetric positive definite $d^{\mathbf{R}} \times d^{\mathbf{R}}$ matrices, which can be accomplished by re-parameterising the $\boldsymbol{\Sigma}$ matrix (Bateman & Pinheiro, 2000[55](#)). For the general case where d -random effects are involved, the following procedure is recommended:

1. Before conducting any optimisation, convert $\boldsymbol{\Sigma}$ to the unconstrained space $\boldsymbol{\theta}$.
When:

(a) $d = 1$:

$$\boldsymbol{\theta} \equiv \frac{1}{2} \log(\boldsymbol{\Sigma}).$$

(b) $d > 1$:

i. Obtain the spectral decomposition of $\boldsymbol{\Sigma}$,

$$\boldsymbol{\Sigma} = \mathbf{u}_{\boldsymbol{\Sigma}} \boldsymbol{\lambda}_{\boldsymbol{\Sigma}} \mathbf{u}_{\boldsymbol{\Sigma}}^{\top},$$

where $\boldsymbol{\lambda}_{\boldsymbol{\Sigma}}$ is the diagonal matrix of eigenvalues and $\mathbf{u}_{\boldsymbol{\Sigma}}$ is the orthogonal matrix of matching eigenvectors.

ii. Calculate the matrix logarithm of $\boldsymbol{\Sigma}$ using the spectral decomposition,

$$\log(\boldsymbol{\Sigma}) = \mathbf{u}_{\boldsymbol{\Sigma}} \log(\boldsymbol{\lambda}_{\boldsymbol{\Sigma}}) \mathbf{u}_{\boldsymbol{\Sigma}}^{\top}.$$

iii. Use the logarithm $\boldsymbol{\Sigma}$ to acquire $\boldsymbol{\theta}$,

$$\boldsymbol{\theta} \equiv \text{vech}\left(\frac{1}{2} \log(\boldsymbol{\Sigma})\right).$$

2. Now use a quasi-Newton optimisation method to obtain the maximum likelihood approximation

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}) = \underset{\boldsymbol{\theta} \in \mathbb{R}^{d(d+1)/2}}{\operatorname{argmax}} \ell(\boldsymbol{\beta}, \boldsymbol{\theta}).$$

We suggest conducting an initial search via the Nelder-Mead method, with refinements by BFGS algorithm. Both can be implemented via the `optim()` R function in the “stats” package (R Core Team, 2019^[56]).

3. Since it is well established that working with $\log(\text{standard deviation})$ and $\tanh^{-1}(\text{correlation})$ parameters for confidence interval construction results in better asymptotic normality than their non-transformed counterparts, convert $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}})$ to $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}})$. For:

(a) $d = 1$, no conversion is required since there is no correlation parameter and the variance parameter is already log transformed. For consistency with the multidimensional setting, note the following notation change:

$$\hat{\boldsymbol{\omega}} = \hat{\boldsymbol{\theta}}.$$

(b) $d > 1$:

i. Obtain the spectral decomposition of $\text{vech}^{-1}(\hat{\boldsymbol{\theta}})$,

$$\text{vech}^{-1}(\hat{\boldsymbol{\theta}}) = \mathbf{u}_{\hat{\boldsymbol{\theta}}} \boldsymbol{\lambda}_{\hat{\boldsymbol{\theta}}} \mathbf{u}_{\hat{\boldsymbol{\theta}}}^{\top},$$

where $\lambda_{\hat{\theta}}$ is the diagonal matrix of eigenvalues and $\mathbf{u}_{\hat{\theta}}$ is the orthogonal matrix of matching eigenvectors.

- ii. Then using the spectral decomposition find exponent of $\lambda_{\hat{\theta}}$,

$$\hat{\Sigma} = \mathbf{u}_{\hat{\theta}} \exp(2\lambda_{\hat{\theta}}) \mathbf{u}_{\hat{\theta}}^{\top}.$$

- iii. Convert to $\hat{\omega}$

$$\hat{\omega} = \begin{bmatrix} \frac{1}{2} \log(\text{diag}(\hat{\Sigma})) \\ \tanh^{-1} \left\{ \text{vecbd}(\hat{\Sigma}) / \sqrt{\text{vecbd}(\text{diag}(\hat{\Sigma}) \text{diag}(\hat{\Sigma})^{\top})} \right\} \end{bmatrix}.$$

4. Find the Hessian matrix at the maximum $(\hat{\beta}, \hat{\omega})$ (denoted by $H\ell(\hat{\beta}, \hat{\omega})$) using the BFGS quasi-Newton method, which returns the hessian when implemented via the R function `optim()` in the “stats” package (R Core Team, 2019⁵⁶). Although we need values of $(\hat{\beta}, \hat{\omega})$ to be returned in the Hessian, the constraints on these parameters mean the Hessian should be calculated on the $(\hat{\beta}, \hat{\theta})$ space. The conversion can be obtained as follows, for:

- (a) $d = 1$:

$$\hat{\theta} = \exp(2\hat{\omega}).$$

- (b) $d > 1$, form the $d \times d$ symmetric matrix $\hat{\Sigma}$:

- i. Let $\hat{\omega}_1$ denote the first d entries of $\hat{\omega}$ and $\hat{\omega}_2$ denote the remaining $\frac{1}{2}d(d-1)$ entries of $\hat{\omega}$. Set $\text{diag}(\hat{\Sigma}) = \exp(2\hat{\omega}_1)$. Obtain the below-diagonal entries of $\hat{\Sigma}$ so that

$$\text{vecbd}(\hat{\Sigma}) = \tanh(\hat{\omega}_2) \odot \text{vecbd}(\exp(\hat{\omega}_1) \exp(\hat{\omega}_1)^{\top})$$

holds. Obtain the above-diagonal entries of $\hat{\Sigma}$, and reconstruct $\hat{\Sigma}$ by symmetry.

- ii. Find the spectral decomposition of $\hat{\Sigma}$,

$$\hat{\Sigma} = \mathbf{u}_{\hat{\Sigma}} \text{diag}(\lambda_{\hat{\Sigma}}) \mathbf{u}_{\hat{\Sigma}}^{\top}.$$

- iii. Using the spectral decomposition make the conversion to θ ,

$$\hat{\theta} = \text{vech} \left\{ \frac{1}{2} \mathbf{u}_{\hat{\Sigma}} \text{diag}(\log(\lambda_{\hat{\Sigma}})) \mathbf{u}_{\hat{\Sigma}}^{\top} \right\}.$$

5. Form $100(1 - \alpha)\%$ confidence intervals for the entries of $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}})$ using

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\omega}} \end{bmatrix} \pm \Phi^{-1} \left(1 - \frac{1}{2} \alpha \right) \sqrt{-\text{diag} \left\{ (\mathbf{H} \ell(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}))^{-1} \right\}}.$$

6. Back transform the confidence interval limits for the $\hat{\boldsymbol{\omega}}$ component, to correspond to the standard deviation and correlation parameters as follows:

$$\left[\frac{\sqrt{\text{diag}(\hat{\boldsymbol{\Sigma}})}}{\text{vecbd}(\hat{\boldsymbol{\Sigma}}) / \sqrt{\text{vecbd}(\text{diag}(\hat{\boldsymbol{\Sigma}})\text{diag}(\hat{\boldsymbol{\Sigma}})^\top)}} \right].$$

3.3 Best predictor

In addition to estimating $(\boldsymbol{\beta}, \boldsymbol{\Sigma})$, we also wish to predict the random effects \mathbf{u}_i for each group. For the binary mixed model via EP, the best predictor of \mathbf{u}_i for the multivariate case is found in a similar manner to the univariate case,

$$\begin{aligned} \text{BP}(\mathbf{u}_i) &= E(\mathbf{u}_i | \mathbf{y}_i) \\ &= \int_{\mathbb{R}^{d\mathbf{R}}} \mathbf{u}_i p(\mathbf{u}_i | \mathbf{y}_i; \boldsymbol{\beta}, \boldsymbol{\Sigma}) d\mathbf{u}_i \\ &= \int_{\mathbb{R}^{d\mathbf{R}}} \mathbf{u}_i \left\{ \frac{p(\mathbf{y}_i | \mathbf{u}_i; \boldsymbol{\beta}) p(\mathbf{u}_i; \boldsymbol{\Sigma})}{\int_{\mathbb{R}^{d\mathbf{R}}} p(\mathbf{y}_i | \mathbf{u}_i; \boldsymbol{\beta}) p(\mathbf{u}_i; \boldsymbol{\Sigma})} \right\} d\mathbf{u}_i, \end{aligned} \quad (3.25)$$

where $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})$. We now show that by products of Algorithm [6](#) facilitate the empirical best predictions of the \mathbf{u}_i . Let

$$\hat{\boldsymbol{\eta}}_i \equiv \boldsymbol{\eta}_\Sigma + \text{SUM}(\boldsymbol{\eta}_{p(\mathbf{y}_i | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}) = \begin{bmatrix} \hat{\boldsymbol{\eta}}_{i1} \\ \hat{\boldsymbol{\eta}}_{i2} \end{bmatrix}, \quad (3.26)$$

where $\hat{\boldsymbol{\eta}}_{i1}$ corresponds to the first $d^{\mathbf{R}}$ entries of $\hat{\boldsymbol{\eta}}_i$, $\hat{\boldsymbol{\eta}}_{i2}$ corresponds to the remaining $d^{\mathbf{R}}$ entries of $\hat{\boldsymbol{\eta}}_i$ and $\boldsymbol{\eta}_\Sigma$ and $\text{SUM}(\boldsymbol{\eta}_{p(\mathbf{y}_i | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i})$ are as previously defined in Algorithm [6](#) with $(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = (\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Sigma}})$. Note that Algorithm [6](#) involves using

$$\exp \left\{ \begin{bmatrix} \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \text{vech}(\mathbf{u}_i \mathbf{u}_i^\top) \end{array} \right]^\top \\ \hat{\boldsymbol{\eta}}_i \end{bmatrix} \right\} \text{ to replace } p(\mathbf{y}_i | \mathbf{u}_i; \boldsymbol{\beta}) p(\mathbf{u}_i; \boldsymbol{\Sigma}).$$

Subsequently as in the univariate case, we can approximate $\text{BP}(\mathbf{u}_i) = E(\hat{\mathbf{u}}_i)$, where $\hat{\mathbf{u}}_i$ is multivariate normal with natural parameter $\hat{\boldsymbol{\eta}}_i$. Thus,

$$\text{BP}(\mathbf{u}_i) = -\frac{1}{2} \left(\text{vec}^{-1} \left((\mathbf{D}_d^+)^{\top} \hat{\boldsymbol{\eta}}_{i2} \right) \right)^{-1} \hat{\boldsymbol{\eta}}_{i1}.$$

Although $\text{Cov}(\text{BP}(\mathbf{u}_i) - \mathbf{u}_i)$ is well approximated with $E_{\mathbf{y}_i}(\text{Cov}(\mathbf{u}_i | \mathbf{y}_i))$ (McCulloch, Searle & Neuhaus, 2008^[41]), where

$$\text{Cov}(\mathbf{u}_i | \mathbf{y}) = -\frac{1}{2} \left(\text{vec}^{-1} \left((\mathbf{D}_d^+)^{\top} \hat{\boldsymbol{\eta}}_{i2} \right) \right)^{-1},$$

the approximation is hindered by the expectation over the distribution over the \mathbf{y}_i vector.

3.4 Simulation study

Three simulation studies were conducted, each for 1000 replicates on two computers:

- Computer 1 - MacBook Air laptop with two 2.2 gigahertz processors and 8 gigabytes of random access memory.
- Computer 2 - University of Technology Sydney Interactive High Performance Computing facility Jupiter node with eight 3.6 gigahertz processors and 32 gigabytes random access memory.

We will refer to these computers by name from here on. For each EP algorithm we set our error tolerance value to 10^{-5} and use a maximum of 100 Nelder-Mead search iterations during the optimisation process.

3.4.1 Comparison of maximum likelihood estimates for univariate random effects

The first simulation study was repeated on datasets simulated according to equation (3.1) with true parameter values,

$$\boldsymbol{\beta}_{\text{true}} = [0, 1]^{\top} \quad \text{and} \quad \boldsymbol{\Sigma}_{\text{true}} = \sigma_{\text{true}}^2 = 1. \quad (3.27)$$

Each dataset had $m = 100$ groups and $n_i = 2$ observations in each group. The $\mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{x}_{ij}^{\mathbf{R}}$ vectors were respectively of the form

$$\mathbf{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}]^{\top} \quad \text{and} \quad \mathbf{x}_{ij}^{\mathbf{R}} = 1,$$

where x_{1ij} was generated independently from a uniform distribution on the unit interval. The model described was fit using each of the following methods:

1. Adaptive Gauss-Hermite quadrature using 100 points of evaluation, implemented via the function `glmer()` in the R package “lme4” (Bates, et al., 2018⁵).
2. Laplace approximation implemented via the function `glmer()` in the R package “lme4” (Bates, et al., 2018⁵).
3. EP as described in this Section.
4. Data cloning as used by the R package “dclone” (Solymos, 2010⁶¹) with 10 clones.

We first compare point estimates and confidence intervals produced by adaptive Gauss-Hermite quadrature against Laplace approximation, EP and data cloning. The first row of Figure 3.3 shows Laplace approximation results in poor statistical inference for all parameters estimated, with empirical coverage less than 95% level. The empirical coverage of the variance parameter by Laplace approximation is particularly low at 81.1%. The second row of Figure 3.3 shows EP produces very similar results to adaptive Gauss-Hermite quadrature for both the empirical coverage and confidence intervals of fixed effects parameters. It also shows EP is slightly conservative with coverage 97.5% for the variance parameter compared to 94.2% that of adaptive Gauss-Hermite quadrature. Data cloning produced similar empirical coverage values to EP (95.7%, 95.1% and 97.4% for β_0 , β_1 and σ respectively). Out of the methods tested, data cloning and EP were the only methods with empirical coverage greater than or equal to 95% across all parameters.

Figure 3.2 compares estimated mean squared error and mean squared error of prediction for the four approaches, where the latter is given by the mean squared error of five randomly selected u_i random intercepts. We also include t-based 95% confidence intervals to provide an indication of the variability of simulation-based mean squared error estimation. The plots show EP performs well in comparison with adaptive Gauss-Hermite quadrature for maximum likelihood estimation and best prediction, and generally improves upon Laplace approximation and data cloning. Although we are aware that a comparison of time is obscured by factors such as language of implementation,

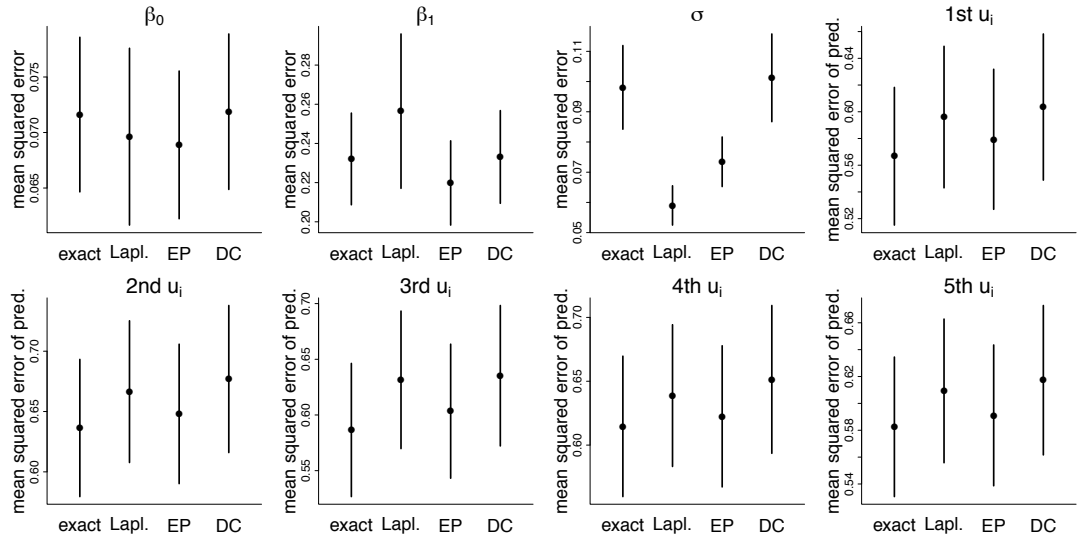


Figure 3.2: Summary of confidence interval coverage for the simulation study with true parameter values from equation (3.29). The horizontal lines are the EP-based confidence intervals for 50 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates and the vertical lines indicate true parameter values. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

the average time for the EP routine over the 1000 replications completed on Computer 1 was 0.1960 seconds, whilst data cloning and adaptive Gauss-Hermite took an average of 143 and 16.1 seconds respectively. Laplace was also extremely time efficient, taking on average 0.158 seconds. It is worth noting construction of confidence intervals occupies a large amount of time for the quadrature approach as it is not optimised with the `glmer()` function.

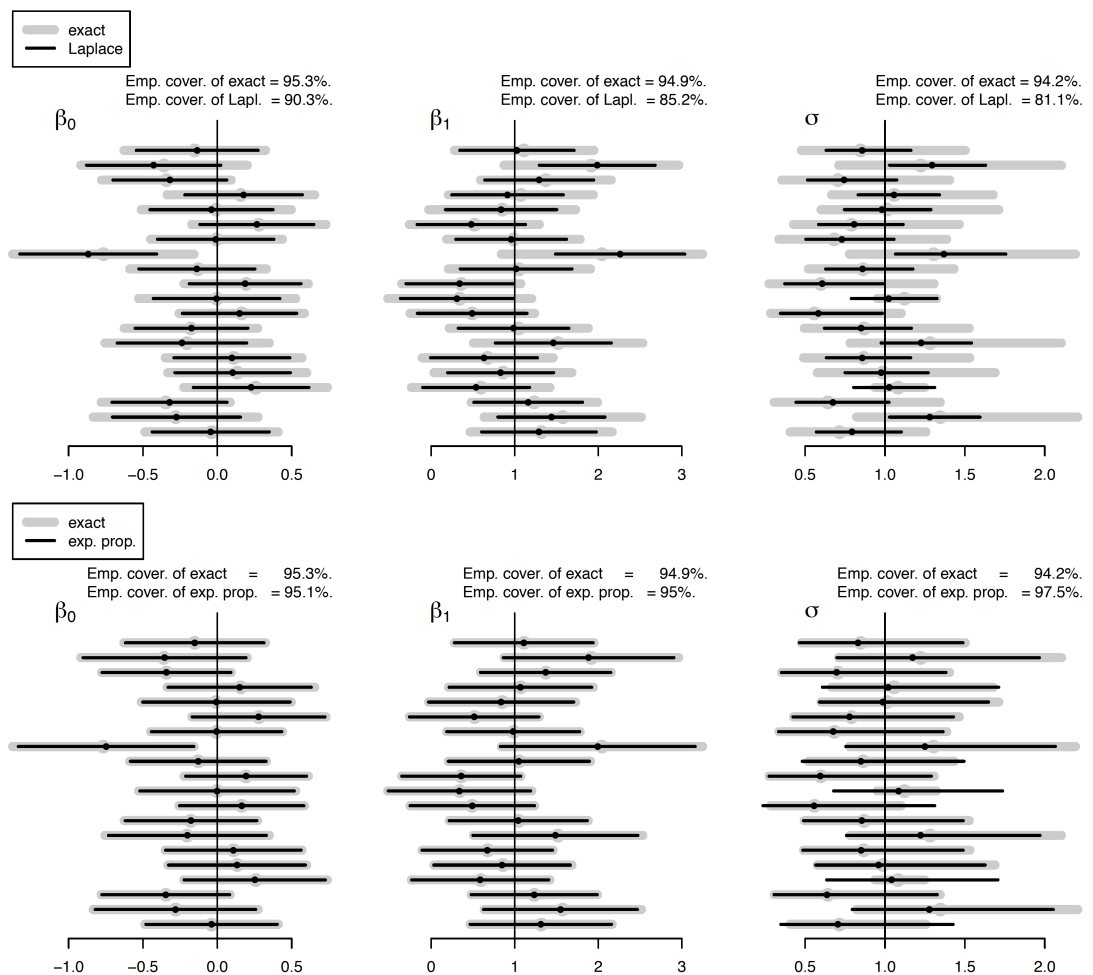


Figure 3.3: Comparison of point estimates and 95% confidence intervals for the simulation study with true parameter values given by equation (3.27). We display 20 randomly chosen replications of the simulation study described. The first row of panels compares exact maximum likelihood with Laplace approximation while the second row compares exact maximum likelihood with expectation propagation approximation. The vertical lines indicate true parameter values and the percentages displayed at the top of each panel are empirical coverages over all 1000 replications for each method involved in the comparison.

3.4.2 Maximum likelihood estimates for bivariate random effects

Data for the second simulation study was simulated according to equation (3.1) with true parameter values

$$\boldsymbol{\beta}_{\text{true}} = [0.37, 0.93, -0.46, 0.08, -1.34, 1.09]^\top \quad \text{and} \quad \boldsymbol{\Sigma}_{\text{true}} = \begin{bmatrix} 0.53 & -0.36 \\ -0.36 & 0.92 \end{bmatrix}. \quad (3.28)$$

We fixed the number of groups in the data at $m = 250$ and the number of measurements in the i th group was a randomly generated integer between 20 and 30 on a uniform distribution. The $\boldsymbol{x}_{ij}^{\mathbf{F}}$ and $\boldsymbol{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\boldsymbol{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}, x_{2ij}, x_{3ij}, x_{4ij}, x_{5ij}]^\top \quad \text{and} \quad \boldsymbol{x}_{ij}^{\mathbf{R}} = [1, x_{1ij}]^\top,$$

where each x_{kij} was generated independently from a uniform distribution on the unit interval.

Figure 3.4 displays point estimates and corresponding 95% confidence intervals for each interpretable model parameters, for 50 randomly chosen replications. As before, the empirical coverage values based on all 1000 replications are in the top right-hand corner of each panel. For five of the nine parameters the empirical coverage values of EP is greater than the expected 95% empirical coverage. However, even the random intercept variance parameter with the lowest coverage of 94.2% is within 1% of the expected coverage. By implementing the required EP algorithm in the low level language Fortran 77 we maintain high speed inference despite the higher samples and increased complexity of the model. On Computer 1, the average computing time over the 1000 replications of the simulation study was 18 seconds, the upper quartile was 20 seconds and the maximum was 34 seconds.

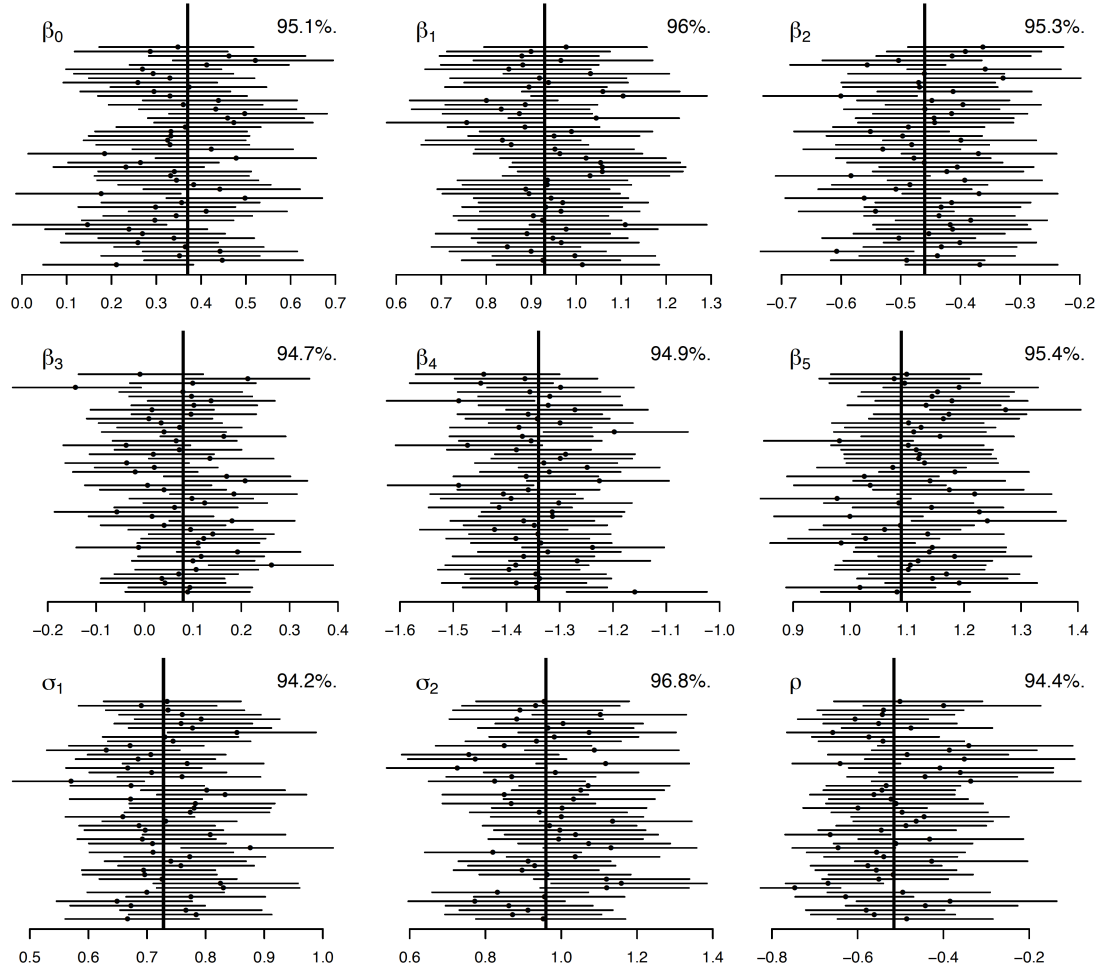


Figure 3.4: Summary of confidence interval coverage for the simulation study with true parameter values from equation (3.28). The horizontal lines are the EP-based confidence intervals for 50 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates and the vertical lines indicate true parameter values. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

3.4.3 Maximum likelihood estimates for trivariate random effects

The third simulation study was repeated 1000 times, where datasets were simulated according to (3.29) with true parameter values

$$\boldsymbol{\beta}_{\text{true}} = [0.37, 0.93, -0.46, 0.08, -1.34, 1.09]^{\top} \quad \text{and} \quad \boldsymbol{\Sigma}_{\text{true}} = \begin{bmatrix} 0.53 & -0.36 & -0.11 \\ -0.36 & 0.92 & 0.08 \\ -0.11 & 0.08 & 0.74 \end{bmatrix}. \quad (3.29)$$

Each dataset had $m = 250$ groups, and each i th group had a randomly generated number of observations, between 20 and 30 on a uniform distribution. The $\boldsymbol{x}_{ij}^{\mathbf{F}}$ and $\boldsymbol{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\boldsymbol{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}, x_{2ij}, x_{3ij}, x_{4ij}, x_{5ij}]^{\top} \quad \text{and} \quad \boldsymbol{x}_{ij}^{\mathbf{R}} = [1, x_{1ij}, x_{2ij}]^{\top},$$

where each x_{kij} was generated independently from a uniform distribution on the unit interval.

The resulting estimates and 95% confidence intervals for each interpretable model parameter of the study are presented in Figure 3.5, where the numbers in the upper-right hand corner of each panel are the empirical coverage values based on all 1000 replicates. Only 25 randomly chosen replicates are shown in each of the panels for ease of viewing. Across all 12 parameters estimated, the empirical coverage values showed excellent accuracy. While three of the fixed effects parameters had less than 95% empirical coverage, three of the random effects parameters had more than 96% empirical coverage. This component of the study was run on Computer 2 and thus we do not provide computational times.

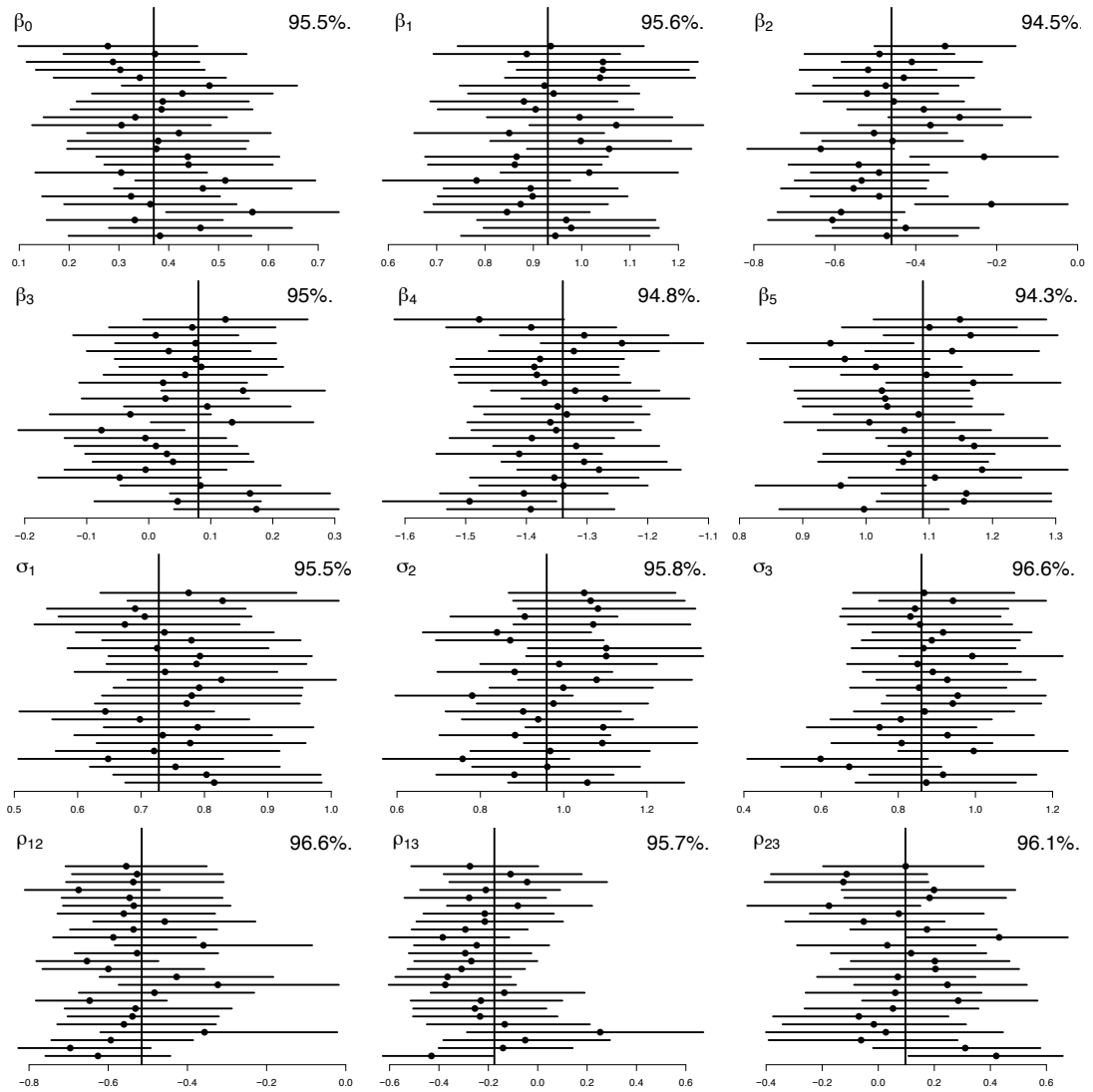


Figure 3.5: Summary of confidence interval coverage for the simulation study with true parameter values from equation (3.29). The horizontal lines are the EP-based confidence intervals for 25 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates and the vertical lines indicate true parameter values. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

3.5 Appendix

3.5.1 Proof of Definition 13

To facilitate the computation of intractable integrals that arise in each of the moments, we now introduce Lemmas 2 and 3. For $\mathbf{x} \in \mathbb{R}^d$, define

$$\phi_{\Sigma}(\mathbf{x}) \equiv (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2} \mathbf{x}^{\top} \Sigma^{-1} \mathbf{x}\right).$$

Lemma 2. For any function $g : \mathbb{R} \rightarrow \mathbb{R}$ and $d \times 1$ vectors α_1 , α_2 and α_3 , the following is true:

$$\int_{\mathbb{R}^d} g(\alpha_1^{\top} \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} g(\|\alpha_1\|x) \phi(x) dx, \quad (3.30)$$

$$\int_{\mathbb{R}^d} g(\alpha_1^{\top} \mathbf{x}) (\alpha_2^{\top} \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = (\alpha_1^{\top} \alpha_2) / \|\alpha_1\| \int_{-\infty}^{\infty} x g(\|\alpha_1\|x) \phi(x) dx \quad (3.31)$$

and

$$\begin{aligned} \int_{\mathbb{R}^d} g(\alpha_1^{\top} \mathbf{x}) (\alpha_2^{\top} \mathbf{x}) (\alpha_3^{\top} \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} &= (\alpha_2^{\top} \alpha_3) \int_{-\infty}^{\infty} g(\|\alpha_1\|x) \phi(x) dx \\ &+ (\alpha_1^{\top} \alpha_2) (\alpha_1^{\top} \alpha_3) / \|\alpha_1\|^2 \int_{-\infty}^{\infty} (x^2 - 1) g(\|\alpha_1\|x) \phi(x) dx. \end{aligned} \quad (3.32)$$

3.5.1.1 Proof of Lemma 2

Note that the integrals on the left-hand side in equations (3.30), (3.31) and (3.32) are respectively

$$E\left(g(\alpha_1^{\top} \mathbf{x})\right), \quad E\left(g(\alpha_1^{\top} \mathbf{x}) (\alpha_2^{\top} \mathbf{x})\right) \quad \text{and} \quad E\left(g(\alpha_1^{\top} \mathbf{x}) (\alpha_2^{\top} \mathbf{x}) (\alpha_3^{\top} \mathbf{x})\right),$$

where

$$\mathbf{x} \sim \mathbf{N}(\mathbf{0}_d, \mathbf{I}_d).$$

We now present simplification of the integral in equation (3.32). Simplification of the integrals in equation (3.30) and (3.31) are analogous. Make the change of variables

$$\mathbf{s} \equiv \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} = \mathbf{A} \mathbf{x}, \quad \text{where} \quad \mathbf{A} \equiv \begin{bmatrix} \alpha_1^{\top} \\ \alpha_2^{\top} \\ \alpha_3^{\top} \end{bmatrix}$$

so that

$$E\left(g(\alpha_1^{\top} \mathbf{x}) (\alpha_2^{\top} \mathbf{x}) (\alpha_3^{\top} \mathbf{x})\right) = E(g(s_1) s_2 s_3), \quad \text{where} \quad \mathbf{s} \sim \mathbf{N}(\mathbf{0}_3, \mathbf{A} \mathbf{A}^{\top}).$$

We then note that,

$$\begin{aligned} E(g(s_1)s_2s_3) &= \int_{-\infty}^{\infty} g(s_1) \left(\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} s_2s_3p(s_2s_3|s_1)ds_2ds_3 \right) p(s_1)ds_1 \\ &= \int_{-\infty}^{\infty} g(s_1) \left(\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\text{Cov}(s_2, s_3|s_1) + E(s_2|s_1)E(s_3|s_1))ds_2ds_3 \right) \\ &\quad \times p(s_1)ds_1 \end{aligned}$$

and

$$\begin{aligned} \begin{bmatrix} s_2 \\ s_3 \end{bmatrix} \Big|_{s_1} &\sim \mathbf{N} \left((s_1/\|\boldsymbol{\alpha}_1\|^2) \begin{bmatrix} \boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_2 \\ \boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_3 \end{bmatrix}, \right. \\ &\quad \left. \begin{bmatrix} \|\boldsymbol{\alpha}_2\|^2 & \boldsymbol{\alpha}_2^\top \boldsymbol{\alpha}_3 \\ \boldsymbol{\alpha}_2^\top \boldsymbol{\alpha}_3 & \|\boldsymbol{\alpha}_3\|^2 \end{bmatrix} - (1/\|\boldsymbol{\alpha}_1\|^2) \begin{bmatrix} (\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_2)^2 & (\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_2)(\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_3) \\ (\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_2)(\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_3) & (\boldsymbol{\alpha}_1^\top \boldsymbol{\alpha}_3)^2 \end{bmatrix} \right). \end{aligned}$$

Simple algebraic manipulations can be used to arrive at equation [3.32](#).

Lemma 3. *For integrals of the forms listed below, corresponding closed form solutions exist:*

$$\int_{\mathbb{R}^d} \Phi(a + \mathbf{b}^\top \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} = \Phi\left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}}\right), \quad (3.33)$$

$$\int_{\mathbb{R}^d} \mathbf{x} \Phi(a + \mathbf{b}^\top \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{b}}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \phi\left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}}\right), \quad (3.34)$$

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^\top \Phi(a + \mathbf{b}^\top \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} \\ = \Phi\left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}}\right) \mathbf{I}_d - \frac{a \mathbf{b} \mathbf{b}^\top}{\sqrt{(\mathbf{b}^\top \mathbf{b} + 1)^3}} \phi\left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}}\right), \end{aligned} \quad (3.35)$$

where $a \in \mathbb{R}$ and \mathbf{b} is $d \times 1$ vector.

3.5.1.2 Proof of Lemma [3](#)

To arrive at equation [\(3.34\)](#), let \mathbf{e}_i be a $d \times 1$ vector with its i th entry equal to 1 and zeroes elsewhere. Then, using Lemma [2](#) the i th entry of the right-hand side of equation [\(3.34\)](#) is

$$\begin{aligned} \int_{\mathbb{R}^d} (\mathbf{e}_i^\top \mathbf{x}) \Phi(a + \mathbf{b}^\top \mathbf{x}) \phi_I(\mathbf{x}) d\mathbf{x} &= \frac{\mathbf{e}_i^\top \mathbf{b}}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} x \Phi(a + \|\mathbf{b}\|x) \phi(x) dx \\ &= -\frac{\mathbf{e}_i^\top \mathbf{b}}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|x) \phi'(x) dx \\ &= \mathbf{e}_i^\top \mathbf{b} \int_{-\infty}^{\infty} \phi(a + \|\mathbf{b}\|x) \phi(x) dx, \end{aligned} \quad (3.36)$$

where equation (3.36) follows via integration by parts. The components of the integrand in equation (3.36) can be expressed as

$$(2\pi)^{-1} \exp \left\{ -\frac{1}{2}((a + \|\mathbf{b}\|x)^2 + x^2) \right\} = \phi \left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \right) \phi \left(\frac{x + a\|\mathbf{b}\|}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \right).$$

Equation (3.34) is of direct consequence.

To arrive at equation (3.35), define \mathbf{e}_j as a $d \times 1$ vector with j th entry equal to 1 and with zeroes elsewhere. Then by equation (3.32), the (i, j) entry of the right-hand side of equation (3.35) is

$$\begin{aligned} & \int_{\mathbb{R}^d} (\mathbf{e}_i^\top \mathbf{x})(\mathbf{e}_j^\top \mathbf{x}) \Phi(a + \mathbf{b}^\top \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} \\ &= (\mathbf{e}_i^\top \mathbf{e}_j) \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|x) \phi(x) dx + \frac{(\mathbf{e}_i^\top \mathbf{b})(\mathbf{e}_j^\top \mathbf{b})}{\|\mathbf{b}\|^2} \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|x) \phi''(x) dx \\ &= (\mathbf{e}_i^\top \mathbf{e}_j) \phi \left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \right) \phi(x) dx - \frac{(\mathbf{e}_i^\top \mathbf{b})(\mathbf{e}_j^\top \mathbf{b})}{\|\mathbf{b}\|^2} \int_{-\infty}^{\infty} \Phi(a + \|\mathbf{b}\|x) \phi'(x) dx, \end{aligned} \quad (3.37)$$

where equation (3.37) follows via integration by parts. The integrand in equation (3.37) is expressible as

$$(2\pi)^{-1} x \exp \left\{ -\frac{1}{2}((a + \|\mathbf{b}\|x)^2 + x^2) \right\} = -x \phi \left(\frac{a}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \right) \phi \left(\frac{x + a\|\mathbf{b}\|}{\sqrt{\mathbf{b}^\top \mathbf{b} + 1}} \right).$$

Equation (3.35) follows. It is easy to prove equation (3.33) using the same method.

Consider the input function from equation (3.6). To obtain the projection of this function to the multivariate normal, we require closed form solutions to the equivalent of the zeroth, first and second moments. Using the $\otimes k$ notation given in equation (1.2) we can calculate each moment as

$$\mathcal{M}_k \equiv \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp((\boldsymbol{\eta}_1^{\text{input}})^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x}) d\mathbf{x} \quad (3.38)$$

and subsequently obtain the optimum natural parameters via algebraic manipulations. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. We now work on the general case for all $\otimes k$. Using the matrix notation defined in equation (1.14)

$$\begin{aligned} \mathcal{M}_k &= \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) (2\pi)^{-d/2} \exp \left\{ \left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array} \right]^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} d\mathbf{x} \\ &\quad \times \exp(A(\boldsymbol{\eta})) (2\pi)^{d/2}. \end{aligned}$$

Using the inverse map of the natural parameters in equation (1.14)

$$\mathcal{M}_k = \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) \phi_{\boldsymbol{\Sigma}}(\mathbf{x} - \boldsymbol{\mu}) d\mathbf{x} (2\pi)^{d/2} \exp(A(\boldsymbol{\eta})).$$

We now implement a change of variable such that $\mathbf{z} = \mathbf{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu})$ thus

$$\mathcal{M}_k = \int_{\mathbb{R}^d} (\boldsymbol{\mu} + \mathbf{\Sigma}^{1/2}\mathbf{z})^{\otimes k} \Phi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + (\mathbf{\Sigma}^{1/2}\mathbf{c}_1)^\top \mathbf{z}) \phi_{\mathbf{I}}(\mathbf{z}) d\mathbf{z} Z_1, \quad (3.39)$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + d/2 \log(2\pi))$. It is then easy to show that each of the k th moments are

$$\begin{aligned} Z_1^{-1} \mathcal{M}_0 &= \Phi(r_2), \\ Z_1^{-1} \mathcal{M}_1 &= \boldsymbol{\mu} \Phi(r_2) + 2\mathbf{\Sigma} \mathbf{c}_1 r_1^{-1} \phi(r_2), \\ Z_1^{-1} \mathcal{M}_2 &= (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \mathbf{\Sigma}) \Phi(r_2) + 2r_1^{-1} (\mathbf{\Sigma} \mathbf{c}_1 \boldsymbol{\mu}^\top + \boldsymbol{\mu} \mathbf{c}_1^\top \mathbf{\Sigma} - 2r_2 r_1^{-1} \mathbf{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \mathbf{\Sigma}) \phi(r_2). \end{aligned}$$

The expressions for the optimal mean and variance parameters follow respectively, as

$$\boldsymbol{\mu}^* = E(\mathbf{x}) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \boldsymbol{\mu} + 2\mathbf{\Sigma} \mathbf{c}_1 r_1^{-1} \zeta'(r_2),$$

and

$$\mathbf{\Sigma}^* = E(\mathbf{x}\mathbf{x}^\top) - E(\mathbf{x})E(\mathbf{x})^\top = \frac{\mathcal{M}_2}{\mathcal{M}_0} - \frac{\mathcal{M}_1}{\mathcal{M}_0} \left(\frac{\mathcal{M}_1}{\mathcal{M}_0} \right)^\top = \mathbf{\Sigma} - (4\mathbf{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \mathbf{\Sigma} r_1^{-2}) \zeta''(r_2).$$

Converting back to non-vectorised natural parameter form using conversion defined in equation (1.7) leads to the required result.

3.5.2 Proof of Result 12

Consider an unnormalised multivariate normal density function $f(\mathbf{x}; \boldsymbol{\eta})$

$$f_{\text{UN}}(\mathbf{x}; \boldsymbol{\eta}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta} \right\}.$$

Then the Kullback-Leibler divergence of f_{UN} from f_{input} is

$$\begin{aligned} \text{KL}(f_{\text{input}} \parallel f_{\text{UN}}) &= \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) \log(f_{\text{input}}(\mathbf{x})/f_{\text{UN}}(\mathbf{x}; \boldsymbol{\eta})) + f_{\text{UN}}(\mathbf{x}; \boldsymbol{\eta}) - f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ &= \mathcal{K}(\boldsymbol{\eta}) + \text{const}, \end{aligned}$$

where

$$\mathcal{K}(\boldsymbol{\eta}) \equiv (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})) - \begin{bmatrix} \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^\top) f_{\text{input}}(\mathbf{x}) d\mathbf{x} \end{bmatrix}^\top \boldsymbol{\eta}.$$

Note that the derivative vector of $\mathcal{K}(\boldsymbol{\eta})$ is

$$\text{DK}(\boldsymbol{\eta}) = (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})) \begin{bmatrix} 1 \\ \text{DA}_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})^\top \end{bmatrix}^\top - \begin{bmatrix} \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^\top) f_{\text{input}}(\mathbf{x}) d\mathbf{x} \end{bmatrix}^\top.$$

Since the stationary condition to minimise the $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$ occurs at $\text{DK}(\boldsymbol{\eta})^\top = \mathbf{0}$,

$$(2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})) \begin{bmatrix} 1 \\ \nabla A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}) \end{bmatrix} = \begin{bmatrix} \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \mathbf{x} f_{\text{input}}(\mathbf{x}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^\top) f_{\text{input}}(\mathbf{x}) d\mathbf{x} \end{bmatrix}, \quad (3.40)$$

where $\nabla A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}) \equiv \text{D}A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})^\top$ is the gradient vector of $A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})$. It is then easy to show

$$\boldsymbol{\eta}_0^* = \log(C_{f_{\text{input}}}) - A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}^*) - \frac{d}{2} \log(2\pi),$$

where

$$\boldsymbol{\eta}_{1:2}^* = (\nabla A_{\mathbf{N}})^{-1} \left(\begin{bmatrix} \int_{\mathbb{R}^d} \mathbf{x} (f_{\text{input}}(\mathbf{x})/C_{f_{\text{input}}}) d\mathbf{x} \\ \int_{\mathbb{R}^d} \text{vech}(\mathbf{x}\mathbf{x}^\top) (f_{\text{input}}(\mathbf{x})/C_{f_{\text{input}}}) d\mathbf{x} \end{bmatrix} \right) \quad (3.41)$$

with existence and uniqueness of $(\nabla A_{\mathbf{N}})^{-1}$ being guaranteed by Proposition 3.2 of Wainwright & Jordan (2008), [66](#) and $C_{f_{\text{input}}} \equiv \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x}$. The Hessian matrix of $\mathcal{K}(\boldsymbol{\eta})$ is

$$\text{HK}(\boldsymbol{\eta}) = (2\pi)^{d/2} \exp(\boldsymbol{\eta}_0 + A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2})) \left\{ \begin{bmatrix} 1 \\ \nabla A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}) \end{bmatrix} \begin{bmatrix} 1 \\ \nabla A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}) \end{bmatrix}^\top + \begin{bmatrix} 0 & \mathbf{0}^\top \\ \mathbf{0} & \text{HA}_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}) \end{bmatrix} \right\}$$

By Proposition 3.1 of Wainwright & Jordan (2008), [66](#) $A_{\mathbf{N}}$ is strictly convex on its domain and thus $\text{HA}(\boldsymbol{\eta}_{1:2})$ is positive definite. Thus, $\text{HK}(\boldsymbol{\eta})$ is positive definite for all $\boldsymbol{\eta}$ and so equation [2.40](#) is the unique maximiser of $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$. Therefore,

$$\text{proj}_{\text{UN}}[f_{\text{input}}](x) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta}^* \right\}$$

where $\boldsymbol{\eta}^*$ is as previously defined. However, $\boldsymbol{\eta}_{1:2}^*$ is the same natural parameter vector that arises via projection of $f_{\text{input}}/C_{f_{\text{input}}}$ onto the family of univariate normal density functions, thus

$$\text{proj}_{\mathbf{N}}[f_{\text{input}}/C_{f_{\text{input}}}] (x) = \exp \left\{ \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix} \boldsymbol{\eta}_{1:2}^* - A_{\mathbf{N}}(\boldsymbol{\eta}_{1:2}^*) \right\} (2\pi)^{-d/2}$$

which immediately leads to Result [3](#).

Chapter 4

Expectation propagation for one level logistic mixed models

Although a probit model can be used to handle data with a binary response variable, logistic models are used more frequently for the elegant interpretation they facilitate. However, since the integral arising in the projection required for fitting logistic GLMMs via EP does not have a closed form solution, the models explored in this section will be more computationally intensive than their probit counterparts. We now explore two methods of obtaining the appropriate update for the logistic link function. The benefit of implementing the message passing approach in the previous two chapters becomes obvious over the next two sections, as the logistic extension becomes easily available with minimal algebraic overheads.

This chapter is broken into two main sections. As in Chapter [2](#), we start with the simplest random intercepts only model in Section [4.1](#). The next Section [4.2](#) follows the work presented in Chapter [3](#) and extends our methodology to a general logistic model for any number of fixed and random effects. Since the previous two chapters explain the message passing frameworks in detail, we provide only the updated details and refer back to the previous Chapters when required.

4.1 The simplest logistic mixed model

As in the probit case, we develop our methodology on the simplest GLMM with random intercepts only, where we consider only the parameter σ^2 . For observed values of

$$y_{ij}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

where $y_{ij} \in \{0, 1\}$, the simplified logistic binary mixed model can be shown to be

$$y_{ij}|u_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\text{expit}(u_i)), \quad u_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2),$$

where $\text{expit}(x)$ is given by equation (1.9) and u_i is a scalar unobserved latent variable. We wish to find the maximiser of $\ell(\sigma^2)$ denoted by $\widehat{\sigma^2}$. The log-likelihood can be expressed as a sum

$$\ell(\sigma^2) = \sum_{i=1}^m \ell_i(\sigma^2),$$

where

$$\ell_i(\sigma^2) \equiv \log \int_{-\infty}^{\infty} \prod_{j=1}^n \text{expit}((2y_{ij} - 1)u_i) (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/2\sigma^2) du_i, \quad (4.1)$$

and the maximum likelihood estimate of σ^2 is given by

$$\widehat{\sigma^2} = \underset{\sigma^2}{\text{argmax}} \ell(\sigma^2).$$

The best predictor is given by

$$\text{BP}(u_i) = \frac{\int_{-\infty}^{\infty} u_i \prod_{j=1}^n \text{expit}((2y_{ij} - 1)u_i) \exp(-u_i^2/2\sigma^2) du_i}{\int_{-\infty}^{\infty} \prod_{j=1}^n \text{expit}((2y_{ij} - 1)u_i) \exp(-u_i^2/2\sigma^2) du_i}.$$

Although we are working with the simplest model, the calculation of the maximum likelihood estimator and best predictors are complicated by the intractable integrals arising in equation (4.1). To calculate the log-likelihood surface of $\ell(\sigma^2)$ we explore an EP approximation scheme and compare it to a traditional quadrature approach.

The following subsections copy the structure of Chapter 2. In Subsection 4.1.1 we provide details of the quadrature approach to estimating the likelihood surface, then explain our novel method using EP in Subsection 4.1.2. We compare the likelihood surface of both methods in Subsection 4.1.3 and discuss best predictor computation in Subsection 4.1.4. Point estimate and confidence interval calculation is conducted

analogously to Section 2.4.1 and thus we do not repeat it.

4.1.1 Traditional quadrature likelihood approximation

Implementation of the adaptive quadrature via the Gauss-Kronrod method follows the same approach as in the probit case, where we utilise the R function `integrate()` included with base R. Although it is possible to directly implement the integrate function for calculation of intractable integrals arising in equation (4.1), doing so leads to results of small absolute value and as such are difficult to store. With this in mind, denoting $y_{ij}^\dagger \equiv 2y_{ij} - 1$, we suggest that for numerical stability each $\ell_i(\sigma^2)$ is calculated as

$$\frac{1}{2} \log(2\pi\sigma^2) + \ell_i(\sigma^2) = h(u_{0i}) + \log \int_{-\infty}^{\infty} \exp(h_i(u) - h_i(u_{0i})) du,$$

where

$$\begin{aligned} h_i(u) &\equiv \sum_{j=1}^n \left\{ y_{ij}^\dagger u - \log \left(1 + \exp(y_{ij}^\dagger u) \right) \right\} - \frac{u^2}{2\sigma^2}, \\ h_i'(u) &= \sum_{j=1}^n \left\{ y_{ij}^\dagger \left(1 + \exp(y_{ij}^\dagger u) \right)^{-1} \right\} - \frac{u}{\sigma^2} \end{aligned}$$

and u_{0i} is the root of h_i' which we recommend finding using a bisection search, where the starting values are selected -1 and 1 to be for the lower and upper bounds respectively.

4.1.2 Expectation propagation likelihood approximation

We now consider an EP approach to the approximate likelihood (denoted by $\ell(\sigma^2)$) which follows that of Section 2.2. The EP approximation is motivated by the minimisation of a KL divergence criterion, which in the logistic case is used to select an unnormalised normal density function to replace each

$$\text{expit}((2y_{ij} - 1)u_i), \quad 1 \leq j \leq n$$

in equation (4.1). Subsequently, the integrand is proportional to a product of univariate normal density functions.

The goal of the EP problem is to find the KL projection of the input function onto the family of normal density functions. For logistic binary GLMMs, EP requires

repeated projection of the form

$$f_{\text{input}}(x) = \text{expit}(c_0 + c_1x) \exp(\eta_1^{\text{input}}x + \eta_2^{\text{input}}x^2) \quad (4.2)$$

onto an unnormalised normal distribution (written in exponential form in equation (2.3)), where $\eta_1^{\text{input}} \in \mathbb{R}$ and $\eta_2^{\text{input}} < 0$, and c_0 , c_1 and x follow from the probit case. As such, we seek $\boldsymbol{\eta}^*$ such that

$$\int_{-\infty}^{\infty} x^k \text{expit}(c_0 + c_1x) \exp(\eta_1^{\text{input}}x + \eta_2^{\text{input}}x^2) dx = \int_{-\infty}^{\infty} x^k \exp \left\{ \left[\begin{array}{c} 1 \\ x \\ x^2 \end{array} \right]^{\top} \boldsymbol{\eta}^* \right\} dx, \quad (4.3)$$

where $\boldsymbol{\eta}^*$ is defined in equation (2.4). As before, obtaining the natural parameters $\boldsymbol{\eta}^*$ for projection onto the unnormalised normal family follows from obtaining the projection onto the normal family. As per Result 3, the optimal natural parameters η_1^* and η_2^* , are given according to the projection of the normalised function $f_{\text{input}}/C_{f_{\text{input}}}$ onto the normal family. We can subsequently use these optimal natural parameters to find the normalising natural parameter η_0^* via Result 4 and thus obtain the projection onto unnormalised normal family.

Thus to obtain the required projection, we first obtain the optimal natural parameters η_1^* and η_2^* to project onto the normal family. However, unlike the probit case no closed form solution exists for the integral arising in the required projection (left hand side of equation (4.3)). As in Nolan & Wand (2017)⁴⁷ we turn to the piecewise approximation of the expit function by Monahan & Stefanski (1989).⁴⁶ The key result can be summarised as,

$$\begin{aligned} \text{expit}(x) &\approx \text{expit}_{\text{MS}}(x), \\ \text{where } \text{expit}_{\text{MS}}(x) &= \sum_{i=1}^8 p_i \Phi(s_i x) \end{aligned} \quad (4.4)$$

and p_i and s_i are fixed constants given in Table 18.4.1 of Monahan & Stefanski (1989).⁴⁶ Thus we can write

$$\sum_{i=1}^8 p_i \int_{-\infty}^{\infty} x^k \Phi(s_i c_0 + s_i c_1 x) \exp(\eta_1^{\text{input}}x + \eta_2^{\text{input}}x^2) dx \approx \int_{-\infty}^{\infty} x^k \exp \left\{ \left[\begin{array}{c} 1 \\ x \\ x^2 \end{array} \right]^{\top} \boldsymbol{\eta}^* \right\} dx$$

and implement Lemma 1 to solve the required integrals. Using simple algebraic manipulations analogous to the simplest probit case shown in equation (2.2) we arrive at Result 16

Result 16. Given f_{input} follows the form of equation (4.2), the projection onto the univariate normal family is given by

$$proj_N[f_{input}] = \exp\left(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*)\right) h(x)$$

where

$$\boldsymbol{\eta}_{-1}^* = k_{logistic}(\boldsymbol{\eta}_{-1}^{input}; c_0, c_1), \quad \boldsymbol{\eta}_{-1}^{input} \equiv \begin{bmatrix} \eta_1^{input} \\ \eta_2^{input} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix},$$

$k_{logistic}\left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1\right)$ is as per Definition 15 and $\mathbf{T}(x)$ and $h(x)$ follow from Section 1.5.2.1

Definition 15. For primary arguments $a_1 \in \mathbb{R}$ and $a_2 < 0$ and auxiliary arguments $c_0, c_1 \in \mathbb{R}$ the function $k_{logistic} : H \rightarrow H$ is given by

$$k_{logistic}\left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1\right) = \begin{bmatrix} r_5(a_1 + c_1 r_3) \\ r_5 a_1 \end{bmatrix},$$

where

$$r_{1i} = \sqrt{2(2 - s_i^2 c_1^2 a_2^{-1})}, \quad r_{2i} = s_i(2c_0 - c_1 a_2^{-1} a_1) r_{1i}^{-1}, \quad r_3 = \frac{2 \sum_{i=1}^8 p_i s_i \phi(r_{2i}) r_{1i}^{-1}}{\sum_{i=1}^8 p_i \Phi(r_{2i})},$$

$$\tilde{r}_3 = 4 \sum_{i=1}^8 \left(\frac{p_i s_i^2 r_{2i}}{\sum_{i=1}^8 p_i \Phi(r_{2i})} \right), \quad r_4 = \frac{1}{2}(r_3^2 + \tilde{r}_3), \quad r_5 = (a_2 + r_4 c_1^2)^{-1} a_2.$$

Using Result 16 we now obtain the normalising natural parameter η_0^* to find the projection onto unnormalised normal family.

4.1.2.1 Projection onto the unnormalised normal family

Recall the moment matching problem described by equation (4.3) and Result 4. Then the normalising factor, can be shown to be

$$C_{f_{\text{input}}} = \int_{\mathbb{R}} f_{\text{input}}(x) dx = (2\pi)^{-1/2} \exp(A(\boldsymbol{\eta}^{\text{input}})) \sum_{i=1}^8 p_i \Phi(r_{2i}),$$

where r_{2i} is given in Definition 10 and $A(\boldsymbol{\eta})$ is defined in Section 1.5. By Result 4

$$\eta_0^* = \log \sum_{i=1}^8 p_i \Phi(r_{2i}) + \frac{1}{4}(\eta_1^*)^2/\eta_2^* - \frac{1}{4}(\eta_1^{\text{input}})^2/\eta_2^{\text{input}} + \frac{1}{2} \log(\eta_2^*/\eta_2^{\text{input}}).$$

Thus to obtain η_0^* we introduce Definition 16.

Definition 16. Consider first primary scalar arguments a_1, a_2, b_1 and b_2 , and auxiliary scalar arguments c_0 and c_1 . The function $c_{\text{logistic}} : H \rightarrow \mathbb{R}$ is given by

$$c_{\text{logistic}} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}; c_0, c_1 \right) \equiv \log \left(\sum_{i=1}^8 p_i \Phi(r_{2i}) \right) + \frac{1}{4} b_1^2/b_2 - \frac{1}{4} a_1^2/a_2 + \frac{1}{2} \log(b_2/a_2),$$

where r_{2i} follow from Definition 15.

In summary, we can obtain the projection of the input function onto the univariate normal family using k_{logistic} , then find projection onto the unnormalised univariate normal family using c_{logistic} . This is formalised in Result 17.

Result 17. For an unnormalised input function of the form of equation (4.2),

$$proj_{UN}[f_{input}] = \exp \left\{ \left[\begin{array}{c} 1 \\ x \\ x^2 \end{array} \right]^\top \left[\begin{array}{c} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{array} \right] \right\},$$

where

$$\left[\begin{array}{c} \eta_1^* \\ \eta_2^* \end{array} \right] = k_{logistic} \left(\left[\begin{array}{c} \eta_1^{input} \\ \eta_2^{input} \end{array} \right]; c_0, c_1 \right)$$

and

$$\eta_0^* = c_{logistic} \left(\left[\begin{array}{c} \eta_1^{input} \\ \eta_2^{input} \end{array} \right], \left[\begin{array}{c} \eta_1^* \\ \eta_2^* \end{array} \right]; c_0, c_1 \right).$$

We now explain how to implement the results shown in this section in a message passing framework.

4.1.2.2 Message passing formulation

As mentioned throughout this section, only minor changes to the message passing algorithm discussed in Section 2.2.2 are required to account for the logistic case. Note the components of the sum of the log-likelihood function $\ell_i(\sigma^2)$ are changed to

$$\ell_i(\sigma^2) = \log \int_{-\infty}^{\infty} \left(\prod_{j=1}^n p(y_{ij}|u_i) \right) p(u_i; \sigma^2) du_i, \quad (4.5)$$

where

$$p(y_{ij}|u_i) \equiv \text{expit}((2y_{ij} - 1)u_i) \quad \text{and} \quad p(u_i; \sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/(2\sigma^2))$$

are the conditional density functions of each response given its random effect and the density function of that random effect respectively.

Note the dependence structure of the product in equation (4.5) matches the probit

model shown in Figure 2.1. As such, we only need to update the messages depending on the conditional density functions of each response given its random effect. Following Section 2.2.2, this involves updating messages $m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i)$ (as in equation (2.10)) to

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}_{\text{UN}}[\text{expit}(c_0 + c_{1_{ij}}u_i) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)]}{\exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)},$$

where $c_0 = 0$ and $c_{1_{ij}} \equiv 2y_{ij} - 1$. Utilising Result 6 leads to equation (2.16), where the linear and quadratic coefficient updates in equation (2.17) are changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} \leftarrow k_{\text{logistic}}(\boldsymbol{\eta}_{1:2}^\otimes; c_0, c_{1_{ij}}) - \boldsymbol{\eta}_{1:2}^\otimes \quad (4.6)$$

and where the constant coefficient update in equation (2.18) is changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 \leftarrow c_{\text{logistic}}(\boldsymbol{\eta}_{1:2}^\otimes, (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_0, c_{1_{ij}}). \quad (4.7)$$

In summary, Algorithm 3 applies to the logistic model, however equations (2.17) and (2.18) are replaced with equations (4.6) and (4.7). Additionally, the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$ now follows equation (4.9).

4.1.2.3 Starting values for the univariate logistic case

The EP message passing algorithm proposed relies on good starting values for convergence. We now derive starting values for $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$ using a Taylor series expansion. Note that

$$\log p(y_{ij}|u_i) = f(a_{ij}) - \log(2), \quad \text{where } a_{ij} \equiv (2y_{ij} - 1)u_i,$$

and

$$\begin{aligned} f(x) &= \text{expit}(x), \\ f'(x) &= \frac{\exp(-x)}{(\exp(-x) + 1)^2}, \\ f''(x) &= \frac{\exp(-x)(1 - \exp(-x))}{(\exp(-x) + 1)^3}. \end{aligned} \quad (4.8)$$

Let \hat{u}_i be a Laplace approximation to u_i . Now consider the following Taylor series expansion of the data dependent component of $\ell(\sigma^2)$

$$\begin{aligned} f(a_{ij}) &= f(\hat{a}_{ij}) + (u_i - \hat{u}_i)(2y_{ij} - 1)f'(\hat{a}_{ij}) + \frac{1}{2}((u_i - \hat{u}_i)(2y_{ij} - 1))^2 f''(\hat{a}_{ij}) + \dots \\ &= \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where $\hat{a}_{ij} \equiv (2y_{ij} - 1)\hat{u}_i$ and

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{a}_{ij}) \\ (2y_{ij} - 1)f'(\hat{a}_{ij}) \\ \frac{1}{2}f''(\hat{a}_{ij}) \end{bmatrix}.$$

As in the probit case the quadratic approximation to $\log p(y_{ij}|u_i)$ based on Taylor expansion about \hat{u}_i is $\log \check{p}(y_{ij}|u_i)$, where $\check{p}(y_{ij}|u_i)$ is analogous to equation (2.22). By following the same logic it is easy to show

$$\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ (2y_{ij} - 1)f'(\hat{a}_{ij}) - f''(\hat{a}_{ij})\hat{u}_i \\ \frac{1}{2}f''(\hat{a}_{ij}) \end{bmatrix}, \quad (4.9)$$

where

$$\eta_0^{\text{start}} = f(\hat{a}_{ij}) - (2y_{ij} - 1)f'(\hat{a}_{ij})\hat{u}_i + \frac{1}{2}(2y_{ij} - 1)f''(\hat{a}_{ij})\hat{u}_i^2.$$

The comments at the end of Section 2.2.3 are also applicable for these starting values.

4.1.3 Evaluation of the estimates

Implementing the quadrature and EP approaches in the R computing environment, we now visually compare the accuracy of our likelihood approximation $\check{\ell}'(\sigma^2)$ to the exact likelihood surface $\ell(\sigma^2)$. Figure 4.1 plots the estimates of the likelihood surface for both methods. The data generated had $m = 30$ groups with $n = 5$ responses per group, and the true variance of $\sigma^2 = 0.19$. Using the quadrature method as exact,

the plot shows that although there are some discrepancies on the tails of the likelihood surface, the approximate method follows the exact likelihood surface near the maximum. Additionally, the true value of σ^2 matches well with the maximum of the likelihood. With this in mind, we now turn to methodology for determination of its maximum with 95% confidence intervals.

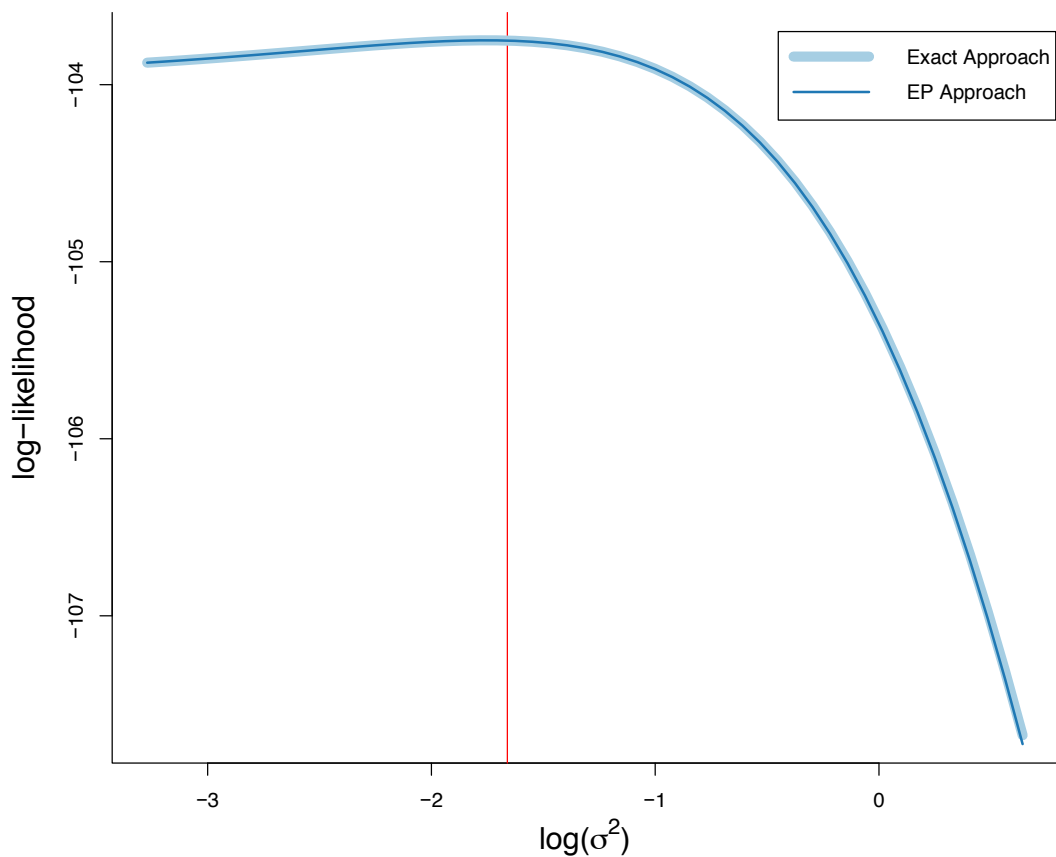


Figure 4.1: A comparison plot of the log-likelihood surface over the parameter σ^2 for logistic models calculated using univariate quadrature and EP via the Monahan & Stefanski (1989)^[46] approximation to the expit function. The true $\sigma^2 = 0.19$ and is represented on the log scale by the red line. The EP approximation is the dark blue line and the exact is the light blue line.

4.1.4 Best predictor

Best prediction of u_i for logistic models via quadrature and EP is the same for the probit model presented in Section 2.5, although for the later we redefine

$$J_{si}(\sigma^2) \equiv \int_{-\infty}^{\infty} u_i^s \prod_{j=1}^n \text{expit}((2y_{ij} - 1)u_i) (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/2\sigma^2) du_i, \quad s = \{0, 1, 2\}.$$

The calculations are otherwise analogous.

4.2 General logistic mixed models

With knowledge from Section 4.1 on the random intercepts only model, we now extend our model to the more general case of GLMMs discussed in Section 1.7 that allow for any number of fixed and random effects. As such we respecify our model as

$$y_{ij} | \mathbf{u}_i \stackrel{\text{ind.}}{\sim} \text{Bernoulli}\left(\text{expit}(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})\right), \quad \mathbf{u}_i \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d\mathbf{R}}, \boldsymbol{\Sigma}), \quad (4.10)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i,$$

where the notation follows the general one level model presented in Section 1.8. The resulting log-likelihood can be written as

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}),$$

where

$$\begin{aligned} \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \log \int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{expit}\left((2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})\right) \right\} \\ &\quad \times |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i, \end{aligned}$$

and the best predictor of \mathbf{u}_i is

$$\text{BP}(\mathbf{u}_i) \equiv \frac{\int_{\mathbb{R}^{d\mathbf{R}}} \mathbf{u}_i \left\{ \prod_{j=1}^{n_i} \text{expit}\left((2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})\right) \right\} \exp\left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i}{\int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{expit}\left((2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})\right) \right\} \exp\left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i} \quad (4.11)$$

Quadrature for integrals greater than two dimensions becomes too computationally intensive for practical implementation, and as such we no longer consider it. The computations required for the multivariate logistic model via EP are mostly the same as the multivariate extension of the probit model, with the only changes being the starting values and projections onto the normalised and unnormalised inverse-logistic density functions.

The structure of this section follows that of the previous, first explaining the schematic of likelihood approximation using EP in Section 4.2.1 and the changes required from the probit model in Section 3.1. It concludes with the results of our simulation studies in Section 4.2.2. Details regarding computation of point estimates and confidence intervals are analogous to those presented in Section 3.2, whilst the same calculations in Section 3.3 can be used for best prediction of \mathbf{u}_i . We do not repeat either section again and instead refer readers to the previous work.

4.2.1 Expectation propagation likelihood approximation

The details of the EP approximation of the likelihood are largely analogous to those provided in Section 3.1 for the probit case. As before, the goal of the EP problem is to find the optimal natural parameters η_0 , $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ which minimise $\text{KL}(f_{\text{input}} \parallel f_{UN})$, where f_{UN} is defined by equation (3.4) and

$$f_{\text{input}}(\mathbf{x}) = \text{expit}(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp\left(\left(\boldsymbol{\eta}_1^{\text{input}}\right)^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x}\right), \quad (4.12)$$

where $\boldsymbol{\eta}_1^{\text{input}}$ is a $d \times 1$ vector and $\mathbf{H}_2^{\text{input}}$ is a $d \times d$ matrix, and c_0 , \mathbf{c}_1 and \mathbf{x} follow from the general probit case. As such, we seek an $\boldsymbol{\eta}^*$ such that

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{expit}(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp\left\{\left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^{\text{input}}\right\} d\mathbf{x} \\ = \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \exp\left\{\left[\begin{array}{c} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^*\right\} d\mathbf{x} \end{aligned} \quad (4.13)$$

where $k \in \{0, 1, 2\}$. To obtain the required projection, we first obtain the optimal natural parameters $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ to project onto the multivariate normal family. For the generalised case, there are two possible solutions worth exploring to solve the otherwise

intractable integrals required for the projections, as per Result [18](#).

Result 18. Given f_{input} follows the form of equation [\(4.12\)](#), the projection onto the multivariate normal family is given by

$$proj_N[f_{input}] = \exp\left(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*)\right) h(\mathbf{x}),$$

where

$$\boldsymbol{\eta}_{-1}^* \equiv K_{logistic_{Quad}}(\boldsymbol{\eta}_{-1}^{input}; c_0, \mathbf{c}_1) \approx K_{logistic_{Approx}}(\boldsymbol{\eta}_{-1}^{input}; c_0, \mathbf{c}_1),$$

$$\boldsymbol{\eta}_{-1}^{input} \equiv \begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix},$$

$K_{logistic_{Quad}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right)$ is defined in Definition [17](#),

$K_{logistic_{Approx}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right)$ is defined in Definition [18](#), and $\mathbf{T}(\mathbf{x})$ and $h(\mathbf{x})$ follow from Section [1.5.2.2](#).

Definition 17. For primary arguments \mathbf{a}_1 ($d \times 1$) and \mathbf{a}_2 ($\frac{1}{2}d(d+1) \times 1$) such that $vec^{-1}\left(-(\mathbf{D}_d^+)^\top \mathbf{a}_2\right)$ is symmetric and positive definite, and auxiliary arguments $c_0 \in \mathbb{R}$ and \mathbf{c}_1 ($d \times 1$), the function $K_{logistic_{Quad}} : H \rightarrow H$ is given by

$$K_{logistic_{Quad}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \equiv \begin{bmatrix} \mathbf{R}_5^\top (\mathbf{a}_1 + r_3 \mathbf{c}_1) \\ \mathbf{D}_d^\top vec(\mathbf{R}_5^\top \mathbf{A}_2) \end{bmatrix}, \quad (4.14)$$

where

$$r_1 = \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{c}_1, \quad r_2 = (2c_0 - \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1) r_1^{-1}, \quad r_3 = r_2 + 2 \mathcal{C}_{b,1:0}(r_6, r_7) r_7,$$

$$r_4 = 2(\mathcal{C}_{b,1:0}(r_6, r_7)^2 - \mathcal{C}_{b,2:0}(r_6, r_7) - r_1) r_1^{-2}, \quad \mathbf{R}_5 = (\mathbf{A}_2 + r_4 \mathbf{c}_1 \mathbf{c}_1^\top)^{-1} \mathbf{A}_2,$$

$$r_6 = 1 - r_2, \quad r_7 = -r_1^{-1}, \quad \mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)},$$

$$\mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, \quad b(x) = \log(1 + e^{-x}),$$

and $\mathcal{C}_b(k, r, q)$ is as per equation [\(4.24\)](#).

Definition 18. For primary arguments \mathbf{a}_1 ($d \times 1$) and \mathbf{a}_2 ($\frac{1}{2}d(d+1) \times 1$) such that $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^{\top} \mathbf{a}_2\right)$ is symmetric and positive definite, and auxiliary arguments $c_0 \in \mathbb{R}$ and \mathbf{c}_1 ($d \times 1$), the function $K_{\text{logistic}_{\text{Approx}}} : H \rightarrow H$ is given by

$$K_{\text{logistic}_{\text{Approx}}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \equiv \begin{bmatrix} \mathbf{R}_5^{\top}(\mathbf{a}_1 + r_3 \mathbf{c}_1) \\ \mathbf{D}_d^{\top} \text{vec}(\mathbf{R}_5^{\top} \mathbf{A}_2) \end{bmatrix}, \quad (4.15)$$

where

$$\begin{aligned} r_{1i} &= \sqrt{2(2 - s_i^2 \mathbf{c}_1^{\top} \mathbf{A}_2^{-1} \mathbf{c}_1)}, & r_{2i} &= s_i(2c_0 - \mathbf{c}_1^{\top} \mathbf{A}_2^{-1} \mathbf{a}_1) r_{1i}^{-1}, \\ r_3 &= \frac{2 \sum_{i=1}^8 p_i s_i r_{1i}^{-1} \phi(r_{2i})}{\sum_{i=1}^8 p_i \Phi(r_{2i})}, & \tilde{r}_3 &= 4 \sum_{i=1}^8 \left(\frac{p_i s_i^2 r_{2i}}{\sum_{i=1}^8 p_i \Phi(r_{2i})} \right), \\ r_4 &= \frac{1}{2}(r_3^2 + \tilde{r}_3), & \mathbf{R}_5 &= (\mathbf{A}_2 + r_4 \mathbf{c}_1^2)^{-1} \mathbf{A}_2. \end{aligned}$$

The first solution using function $K_{\text{logistic}_{\text{Quad}}}$ as per Definition 17 requires solving each of the projections using univariate quadrature. Although this causes computational difficulties, it is still significantly faster than using multivariate quadrature. This approach involves utilising Lemma 2 to express the multivariate integrals required as univariate integrals. To account for the numerically unstable integrals, we express them using the form of $\mathcal{C}_b(p, q, r)$ as in Section 2.1 of Kim & Wand (2018).³²

While optimising the likelihood surface via the EP algorithm, it is not unreasonable for the required integral to be calculated millions of times. As such, speed of integration for each projection becomes a vital issue. For univariate models, adaptive trapezoidal quadrature is a feasible option and provides high accuracy as well as an error criterion. However due to speed limitations its implementation is not practical for multivariate models. Although Gauss-Hermite quadrature with stored weights solves the computational speed issues, it is difficult to implement an error criterion to ensure a reasonable approximation of each integral is obtained. We leave this as an open problem and implement 100 point Gauss-Hermite quadrature for the simulations in this thesis. Details of the algebra required to arrive at Definition 17 are provided in Appendix 4.3.2.

The second approach denoted by $K_{\text{logistic}_{\text{Approx}}}$ as per Definition 18 involves using the piecewise approximation of the expit function by Monahan & Stefanski (1989)⁴⁶ in a similar manner to the previous random intercepts only section. By implementing this approximation, we are able to obtain closed form solutions to the integrals required. Details of the algebra required to arrive at Definition 18 are provided in Appendix 4.3.3.

Practical implementation of this approach in the higher dimensional setting proved to be troublesome with convergence issues arising in the EP algorithm. Figure 18.4.1 of Monahan & Stefanski (1989)^[46] demonstrates the behaviour of the approximation compared to the expit function. An analysis showed that for some values arising in the required projections it is possible for the cumulative normal density function in the expit approximation to become very small, leading to underflow and causing convergence issues in the EP algorithm.

Using Result [18] we now show two methods to obtain the normalising natural parameter and subsequently, the projection onto unnormalised normal family.

4.2.1.1 Projection onto the unnormalised multivariate normal family

Obtaining the projection onto the unnormalised multivariate normal family involves similar algebra to the multivariate probit case. Recall the moment matching problem from equation (4.13) and Results [12] and [13]. We require

$$C_f = \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} = \mathcal{M}_0(c_0, \mathbf{c}_1; \boldsymbol{\eta}),$$

where

$$\mathcal{M}_0(c_0, \mathbf{c}_1; \boldsymbol{\eta}) \equiv \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{expit}(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp(\boldsymbol{\eta}_1^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2 \mathbf{x}) d\mathbf{x},$$

and the right hand side is available by either quadrature or a piecewise approximation.

Using the quadrature method, we express

$$\mathcal{M}_0(c_0, \mathbf{c}_1; \boldsymbol{\eta}) = \mathcal{C}_b(0, r_6, r_7) Z_0 Z_1,$$

where $Z_0 = \exp((r_2/2)^2 r_1 + 0.5 \log(r_7/\pi))$, $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi))$ and r_6 and r_7 are given in Definition [17]. Analogous to the argument given in the univariate case we then get

$$\begin{aligned} \eta_0^* &= \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^\top (\mathbf{H}_2^*)^{-1} \boldsymbol{\eta}_1^* \\ &\quad - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^\top (\mathbf{H}_2^{\text{input}})^{-1} \boldsymbol{\eta}_1^{\text{input}} + \frac{1}{2} \log(|\mathbf{H}_2^*|/|\mathbf{H}_2^{\text{input}}|). \end{aligned}$$

To calculate η_0^* using the quadrature method we introduce Definition [19].

Definition 19. Consider first, primary arguments \mathbf{a}_1 and \mathbf{b}_1 and auxiliary argument \mathbf{c}_1 where all three are $d \times 1$. Next consider arguments \mathbf{a}_2 and \mathbf{b}_2 which are $(\frac{1}{2}d(d+1) \times 1)$ such that both $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^{\top} \mathbf{a}_2\right)$ and $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^{\top} \mathbf{b}_2\right)$ are symmetric and positive definite. Finally note auxiliary argument $c_0 \in \mathbb{R}$. Then the function $C_{\text{logisticQuad}} : H \times H \rightarrow \mathbb{R}$ is given by

$$C_{\text{logisticQuad}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \equiv \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) \\ + \frac{1}{4} \mathbf{b}_1^{\top} \mathbf{B}_2^{-1} \mathbf{b}_1 - \frac{1}{4} \mathbf{a}_1^{\top} \mathbf{A}_2^{-1} \mathbf{a}_1 + \frac{1}{2} \log(|\mathbf{B}_2|/|\mathbf{A}_2|),$$

where $\mathbf{A}_2 \equiv \text{vec}^{-1}\left((\mathbf{D}_d^+)^{\top} \mathbf{a}_2\right)$, $\mathbf{B}_2 \equiv \text{vec}^{-1}\left((\mathbf{D}_d^+)^{\top} \mathbf{b}_2\right)$, r_1 , r_2 , r_6 and r_7 follow from Definition [17](#).

Using the alternative using piecewise approximation,

$$\mathcal{M}_0 = \sum_{i=1}^8 p_i \Phi(r_{2i}) Z_1$$

where Z_1 is defined as in the quadrature approach and r_{2i} is given in Definition [18](#). It is then easy to show

$$\boldsymbol{\eta}_0^* = \log\left(\sum_{i=1}^8 p_i \Phi(r_{2i})\right) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^{\top} (\mathbf{H}_2^*)^{-1} \boldsymbol{\eta}_1^* - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^{\top} (\mathbf{H}_2^{\text{input}})^{-1} \boldsymbol{\eta}_1^{\text{input}} \\ + \frac{1}{2} \log(|\mathbf{H}_2^*|/|\mathbf{H}_2^{\text{input}}|)$$

and arrive at Definition [20](#).

Definition 20. Consider first, primary arguments \mathbf{a}_1 and \mathbf{b}_1 and auxiliary argument \mathbf{c}_1 where all three are $d \times 1$. Next consider arguments \mathbf{a}_2 and \mathbf{b}_2 which are all $(\frac{1}{2}d(d+1) \times 1)$ such that both $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^{\top} \mathbf{a}_2\right)$ and $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^{\top} \mathbf{b}_2\right)$ are symmetric and positive definite. Finally note auxiliary argument $c_0 \in \mathbb{R}$. Then the function $C_{\text{logistic}_{\text{Approx}}}$: $H \times H \rightarrow \mathbb{R}$ is given by

$$C_{\text{logistic}_{\text{Approx}}}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \equiv \log\left(\sum_{i=1}^8 p_i \Phi(r_{2i})\right) + \frac{1}{4} \mathbf{b}_1^{\top} \mathbf{B}_2^{-1} \mathbf{b}_1 - \frac{1}{4} \mathbf{a}_1^{\top} \mathbf{A}_2^{-1} \mathbf{a}_1 + \frac{1}{2} \log(|\mathbf{B}_2|/|\mathbf{A}_2|),$$

where $\mathbf{A}_2 \equiv \text{vec}^{-1}\left((\mathbf{D}_d^+)^{\top} \mathbf{a}_2\right)$, $\mathbf{B}_2 \equiv \text{vec}^{-1}\left((\mathbf{D}_d^+)^{\top} \mathbf{b}_2\right)$, r_1, r_2, r_6 and r_7 follow from Definition [18](#).

For either method, the projection onto the unnormalised multivariate normal family can be obtained with Result [19](#).

Result 19. For an unnormalised input function f_{input} of the form of equation [\(4.12\)](#),

$$\text{proj}_{\text{UN}}[f_{\text{input}}](\mathbf{x}) = \exp\left\{\begin{bmatrix} 1 \\ \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^{\top}) \end{bmatrix}^{\top} \begin{bmatrix} \eta_0^* \\ \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}\right\},$$

where

$$\begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} = K_{\text{logistic}_{\text{Quad}}}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right) \approx K_{\text{logistic}_{\text{Approx}}}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}; c_0, \mathbf{c}_1\right)$$

and

$$\eta_0^* = C_{\text{logistic}_{\text{Quad}}}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1\right) \approx C_{\text{logistic}_{\text{Approx}}}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1\right).$$

4.2.1.2 Message passing formulation

Only minor changes to Algorithm [6](#) are required to account for the logistic case. For the logistic model, the components of the sum of the log-likelihood function $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ follow

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \log \int_{\mathbb{R}^{d_{\mathbf{R}}}} \left(\prod_{j=1}^{n_i} p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) d\mathbf{u}_i, \quad (4.16)$$

where

$$p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \equiv \text{expit} \left((2y_{ij} - 1) (\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}) \right)$$

and

$$p(\mathbf{u}_i; \boldsymbol{\Sigma}) \equiv |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right)$$

are the conditional density functions of each response given its random effect and the density function of that random effect.

As the dependence structure of the product in equation [\(4.16\)](#) is the same as the probit model shown in Figure [3.1](#), we only need to update the messages depending on the conditional density functions of each response given its random effect. This involves updating messages $m_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i)$ (equation [\(3.12\)](#))

$$m_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}_{\text{UN}} \left[\text{expit}(c_{0ij} + \mathbf{c}_{1ij}^\top \mathbf{u}_i) \exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\} \right]}{\exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\}},$$

where $c_{0ij} = (2y_{ij} - 1) \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{c}_{1ij} = (2y_{ij} - 1) \mathbf{x}_{ij}^{\mathbf{R}}$. Following Section [3.1.2](#), we see that this task reduces down to adjusting the calculation of the optimal natural parameters $\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i}$ in equation [\(3.18\)](#), where the linear and quadratic coefficient updates in equation [\(3.19\)](#) given by K_{probit} are changed to either

$$\left(\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{1:2} \leftarrow K_{\text{logistic}_{\text{Quad}}} \left(\boldsymbol{\eta}_{1:2}^\otimes; c_0, \mathbf{c}_{1ij} \right) - \boldsymbol{\eta}_{1:2}^\otimes$$

or

$$\left(\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{1:2} \leftarrow K_{\text{logistic}_{\text{Approx}}} \left(\boldsymbol{\eta}_{1:2}^\otimes; c_0, \mathbf{c}_{1ij} \right) - \boldsymbol{\eta}_{1:2}^\otimes$$

(as per Definitions [17](#) and [18](#) respectively), and where the constant coefficient update in equation [\(3.20\)](#) given by C_{probit} is changed to either

$$\left(\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_0 \leftarrow C_{\text{logistic}_{\text{Quad}}} \left(\boldsymbol{\eta}_{1:2}^\otimes, \left(\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \rightarrow \mathbf{u}_i} \right)_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_0, \mathbf{c}_{1ij} \right)$$

or

$$\left(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}\right)_0 \leftarrow C_{\text{logistic_Approx}} \left(\boldsymbol{\eta}_{1:2}^{\otimes}, \left(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}\right)_{1:2} + \boldsymbol{\eta}_{1:2}^{\otimes}; c_0, \mathbf{c}_{1_{ij}} \right)$$

(as per Definitions [19](#) and [20](#) respectively).

Barring the aforementioned changes and the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ which now follows equation [\(4.17\)](#), the full algorithm for the approximation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ follows Algorithm [6](#).

4.2.1.3 Starting values for the multivariate logistic case

Regardless of the projection methods, the EP algorithm for the logistic model uses the same starting values. We derive starting values using the same principles as in the probit case. Let

$$p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) = f(a_{ij}), \quad \text{where} \quad a_{ij} = (2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}})$$

and $f(x)$ and its derivatives are as defined in equation [\(4.8\)](#). For a Laplace approximation of \mathbf{u}_i denoted by $\hat{\mathbf{u}}_i$ let

$$\hat{a}_{ij} = (2y_{ij} - 1)(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \hat{\mathbf{u}}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}).$$

A Taylor series expansion of $f(a_{ij})$ evaluated at \hat{a}_{ij} leads to

$$\begin{aligned} f(a_{ij}) &= f(\hat{a}_{ij}) + (2y_{ij} - 1)f'(\hat{a}_{ij})(\mathbf{u}_i^\top - \hat{\mathbf{u}}_i^\top)\mathbf{x}_{ij}^{\mathbf{R}} + \frac{1}{2}f''(\hat{a}_{ij})\left((\mathbf{u}_i^\top - \hat{\mathbf{u}}_i^\top)\mathbf{x}_{ij}^{\mathbf{R}}\right)^2 + \dots \\ &= \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}\left((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top\right) \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{a}_{ij}) \\ (2y_{ij} - 1)f'(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}} \\ \frac{1}{2}f''(\hat{a}_{ij})\mathbf{D}_{d_{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix}.$$

Following the same logic of the probit case, quadratic approximation to $\log p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ based on Taylor expansion about $\hat{\mathbf{u}}_i$ is $\log \check{p}(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})$ as per equation [\(3.22\)](#). It is easy

to show

$$\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ (2y_{ij} - 1)f'(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}} - f''(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^{\top}\hat{\mathbf{u}}_i \\ \frac{1}{2}f''(\hat{a}_{ij})\mathbf{D}_{d^{\mathbf{R}}}^{\top}\text{vec}\left(\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^{\top}\right) \end{bmatrix}, \quad (4.17)$$

where

$$\eta_0^{\text{start}} = f(\hat{a}_{ij}) - (2y_{ij} - 1)f'(\hat{a}_{ij})\mathbf{x}_{ij}^{\mathbf{R}}\hat{\mathbf{u}}_i^{\top} + \frac{1}{2}f''(\hat{a}_{ij})\hat{\mathbf{u}}_i^{\top}\mathbf{x}_{ij}^{\mathbf{R}}(\mathbf{x}_{ij}^{\mathbf{R}})^{\top}\hat{\mathbf{u}}_i.$$

The final comments of Section [3.1.3](#) regarding obtaining Laplace approximations to \mathbf{u}_i are also applicable for these starting values.

4.2.2 Simulation study

We now provide a simulation study analogous to that provided for the probit model, where we first compare a quadrature approach to EP and Laplace approximation approach for a random intercept model (i.e. $d^{\mathbf{R}} = 1$) and then test the empirical coverage of the EP approach on a random intercept model and slope model (i.e. $d^{\mathbf{R}} = 2$). For the simulation studies in this chapter we used the quadrature approach to obtaining the projections required for EP and did not assess the speed component.

4.2.2.1 Comparison of maximum likelihood estimates for univariate random effects

The first simulation study was repeated 1000 times, where datasets were simulated with true parameter values:

$$\boldsymbol{\beta}_{\text{true}} = [0.37, 0.93]^{\top} \quad \text{and} \quad \sigma_{\text{true}}^2 = -0.53. \quad (4.18)$$

There were 100 groups generated in the data with each group containing 10 measurements (i.e. $m = 100$, $n = 10$). The $\mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\mathbf{x}_{ij}^{\mathbf{F}} = [1, x_{1,ij}]^{\top} \quad \text{and} \quad \mathbf{x}_{ij}^{\mathbf{R}} = 1$$

where $x_{k,ij}$ was generated independently from a uniform distribution on the unit interval. The tolerance of error values were set to 10^{-5} for the EP scheme and the maximum number of iterations for optimisation was set to 1000. We compared our EP approach

to Laplace approximation and 100 point adaptive Gauss-Hermite quadrature. Both alternative approaches were implemented via the R function `glmer()` from the R package “lme4” (Bates, et al., 2018⁵).

The resulting estimates and 95% confidence intervals for each interpretable model parameter of the study are presented in Figure 4.3, where the numbers in the upper-right hand corner of each panel are the empirical coverage values based on all 1000 replicates. Only 20 randomly chosen replicates from each method are shown in the panels for ease of viewing, where Laplace approximation and EP are shown in black, super imposed adaptive Gauss-Hermite quadrature shown in grey. For the fixed slopes all three methods had the same empirical coverage of 96.3%. The methods also produced similar results for the fixed intercepts, although adaptive Gauss-Hermite quadrature performed the best, with 95.8%, ahead of the 95.7% of EP and 95.4% of Laplace approximation. However, for the variance parameter EP had 97.1% empirical coverage, notably higher Laplace approximation and adaptive Gauss-Hermite quadrature which produced empirical coverages of 93.8% and 94.5% respectively. It appears the results of the three methods are similar for settings of the data, however our method manages to marginally out perform the alternatives based on empirical coverage.

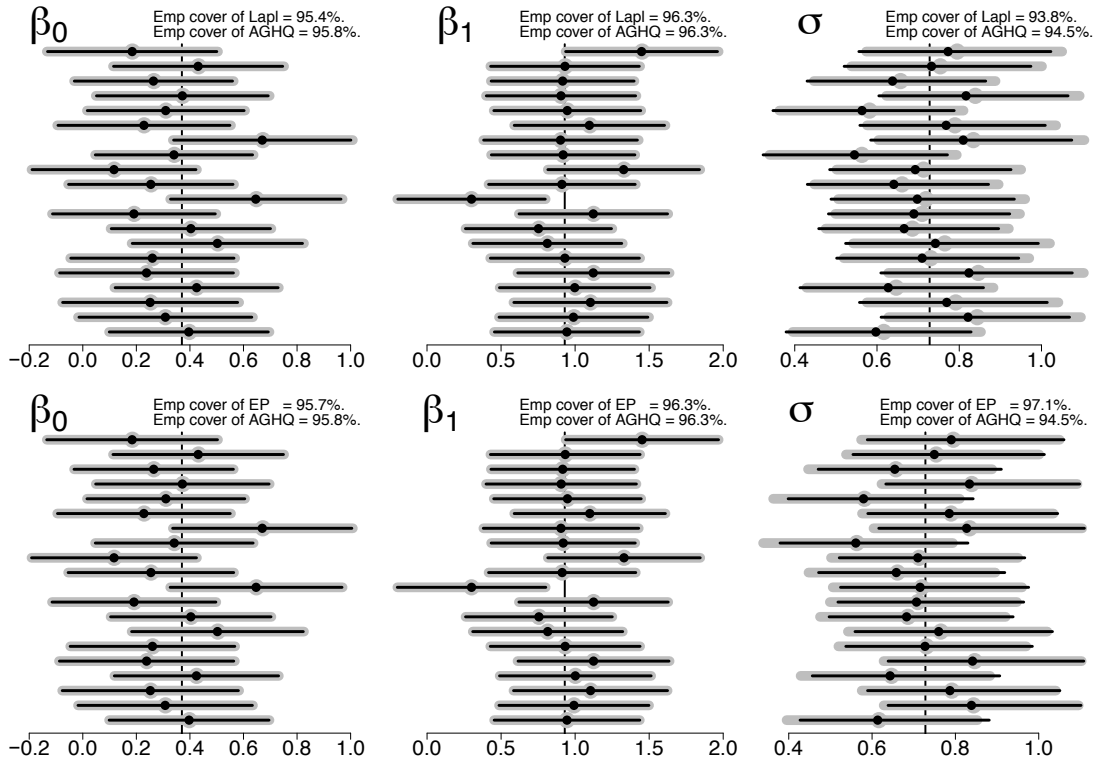


Figure 4.2: Summary comparison between confidence interval coverage for the univariate model with true parameter values from equation (4.18) for models fit with adaptive Gauss-Hermite quadrature and Laplace approximation. The horizontal lines are the confidence intervals for 50 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates, the vertical lines indicate true parameter values, the grey lines correspond to adaptive Gauss-Hermite quadrature, the black lines on the top row correspond to Laplace approximations and the black lines on the bottom row correspond to EP. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

4.2.2.2 Maximum likelihood estimates for bivariate random effects

The simulation study was repeated 1000 times, where datasets were simulated according to equation (4.10) with arbitrarily chosen true parameter values

$$\beta_{\text{true}} = [0.37, 0.93, -0.46, 0.08, -1.34, 1.09]^T \quad \text{and} \quad \Sigma_{\text{true}} = \begin{bmatrix} 0.73 & -0.52 \\ -0.52 & 0.95 \end{bmatrix}. \quad (4.19)$$

We note positive and negative correlation parameters are equally likely, but that neither has a large effect on the model from a mathematical perspective. Additionally, we note that correlations close to 1 pose numerical issues, and in reality we recommend

careful variable selection to avoid them. The number of groups in the data were fixed at $m = 1000$ with six measurements per group. The $\mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\mathbf{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}, x_{2ij}, x_{3ij}, x_{4ij}, x_{5ij}]^{\top} \quad \text{and} \quad \mathbf{x}_{ij}^{\mathbf{R}} = [1, x_{1ij}, x_{2ij}]^{\top}$$

where x_{kij} was generated independently from a uniform distribution on the unit interval. The tolerance of error values were set to 10^{-5} for the EP scheme and the maximum number of iterations for optimisation was set to 1000. The resulting estimates and 95% confidence intervals for each interpretable model parameter of the study are presented in Figure [4.3](#), where the numbers in the upper-right hand corner of each panel are the empirical coverage values based on all 1000 replicates. Only 50 randomly chosen replicates are shown in the each of the panels for ease of viewing. Across all 12 parameters estimated the empirical coverage showed excellent accuracy. While the empirical coverage for 94.3% and 94.6% was below 95% the estimates for the other fixed effects and variance parameters were as promised. The correlation parameter had 96.3% coverage which is slightly conservative.

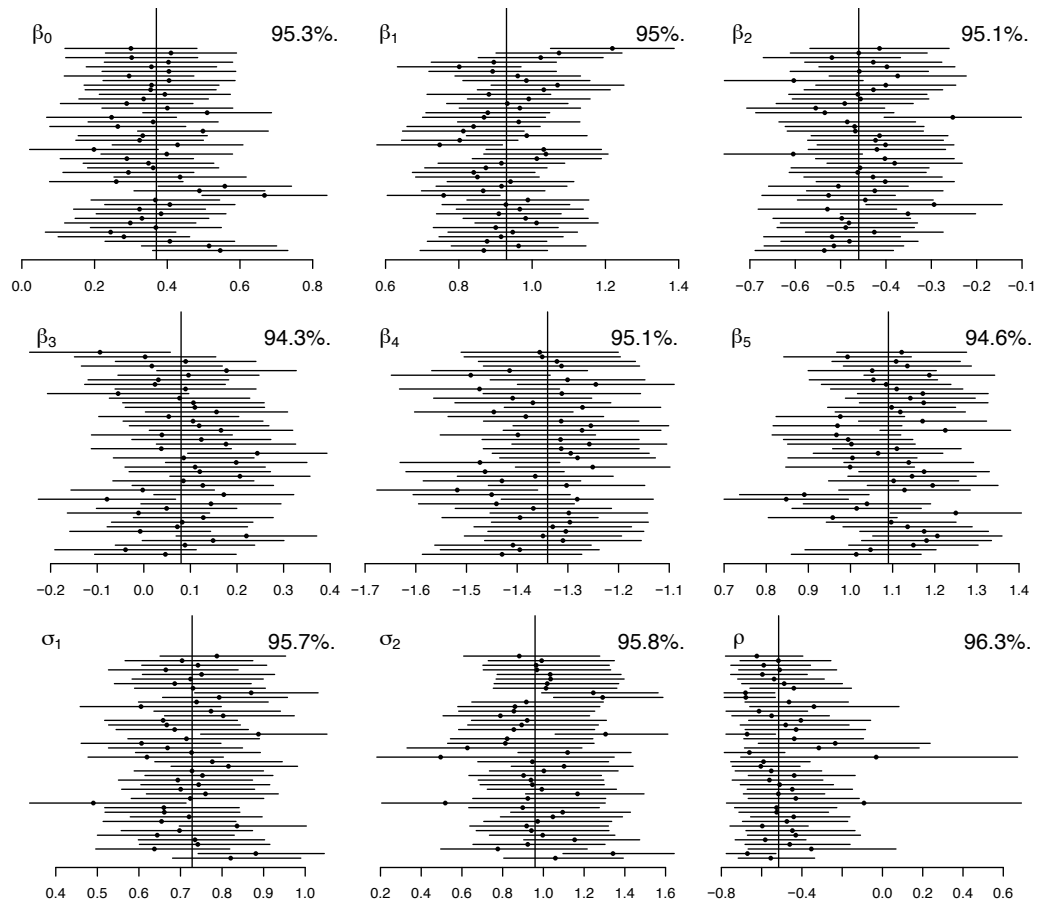


Figure 4.3: Summary of confidence interval coverage for the simulation study with true parameter values from equation (4.19). The horizontal lines are the EP-based confidence intervals for 50 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates and the vertical lines indicate true parameter values. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

4.3 Appendix

4.3.1 Proof of Result [16](#)

Using matrix notation and simple algebraic manipulations we can rewrite the target density of equation [\(4.2\)](#) as

$$\mathcal{M}_k = \sum_{i=1}^8 p_i \int_{-\infty}^{\infty} x^k \Phi(s_i c_0 + s_i c_1 x) \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} dx \exp(A(\boldsymbol{\eta})).$$

Taking the inverse of the natural parameters as per equation [\(1.5\)](#) and implementing a variable change where $u = (x - \mu)\sigma^{-1}$, such that $x = \mu + \sigma u$ and $dx = \sigma du$,

$$\mathcal{M}_k = \sum_{i=1}^8 p_i \int_{-\infty}^{\infty} (\mu + \sigma u)^k \Phi(s_i(c_0 + c_1\mu) + s_i c_1 \sigma u) \phi(u) du Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}))(2\pi)^{1/2}\sigma$. Thus, using Lemma [1](#),

$$\begin{aligned} Z_1^{-1} \mathcal{M}_0 &= \sum_{i=1}^8 p_i \Phi(r_{2i}), \\ Z_1^{-1} \mathcal{M}_1 &= \sum_{i=1}^8 p_i (\mu \Phi(r_{2i}) + 2s_i c_1 \sigma^2 r_{1i}^{-1} \phi(r_{2i})), \\ Z_1^{-1} \mathcal{M}_2 &= \sum_{i=1}^8 p_i (\mu^2 \Phi(r_{2i}) + 4s_i c_1 \mu \sigma^2 r_{1i}^{-1} \phi(r_{2i}) + \sigma^2 \Phi(r_{2i}) - 4s_i^2 c_1^2 \sigma^4 r_{2i} r_{1i}^{-2} \phi(r_{2i})), \end{aligned}$$

where $r_{2i} = 2s_i(c_0 + c_1\mu)r_{1i}^{-1}$ and $r_{1i} = 2((s_i c_1 \sigma)^2 + 1)^{-1/2}$. We can then show that the optimal mean and variance parameters are respectively

$$\mu^* = E(x) = \mu + c_1 \sigma^2 r_3$$

and

$$(\sigma^*)^2 = E(x^2) - (E(x))^2 = \sigma^2 - 2c_1^2 \sigma^4 r_4,$$

where

$$\begin{aligned} r_3 &= \frac{2 \sum_{i=1}^8 s_i p_i \phi(r_{2i}) r_{1i}^{-1}}{\sum_{i=1}^8 p_i \Phi(r_{2i})}, \\ r_4 &= \frac{1}{2} (\tilde{r}_3 + r_3^2) \quad \text{and} \quad \tilde{r}_3 = \frac{4 \sum_{i=1}^8 s_i^2 p_i r_{2i} \phi(r_{2i}) r_{1i}^{-2}}{\sum_{i=1}^8 p_i \Phi(r_{2i})}. \end{aligned}$$

Using matrix notation, converting back to natural parameter form and implementing additional simplification, we arrive at the required result.

4.3.2 Proof of Definition 17

Using simple algebraic manipulations based on Lemma 2, we arrive at Lemma 4:

Lemma 4. *For integrals of the forms listed below, the corresponding solutions exist:*

$$\int_{\mathbb{R}^d} \text{expit}(a + \mathbf{b}^\top \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} \text{expit}(a + \|\mathbf{b}\|x) \phi(x) dx, \quad (4.20)$$

$$\int_{\mathbb{R}^d} \mathbf{x} \text{expit}(a + \mathbf{b}^\top \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{b}}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} x \text{expit}(a + \|\mathbf{b}\|x) \phi(x) dx, \quad (4.21)$$

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^\top \text{expit}(a + \mathbf{b}^\top \mathbf{x}) \phi_{\mathbf{I}}(\mathbf{x}) d\mathbf{x} &= I_d \int_{-\infty}^{\infty} \text{expit}(a + \|\mathbf{b}\|x) \phi(x) dx \\ &+ \frac{\mathbf{b} \mathbf{b}^\top}{\mathbf{b}^\top \mathbf{b}} \left(\int_{-\infty}^{\infty} x^2 \text{expit}(a + \|\mathbf{b}\|x) \phi(x) dx - \int_{-\infty}^{\infty} \text{expit}(a + \|\mathbf{b}\|x) \phi(x) dx \right), \end{aligned} \quad (4.22)$$

where $a \in \mathbb{R}$ and \mathbf{b} is a $d \times 1$ vector.

We wish to obtain the projection of an input function following the form of equation (4.12) onto the multivariate normal family. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. For the general case for all $\otimes k$

$$\mathcal{M}_k \equiv \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{expit}(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp\left(\boldsymbol{\eta}_1^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2 \mathbf{x}\right) d\mathbf{x}. \quad (4.23)$$

Using the matrix notation defined in equation (1.14) we can write

$$\begin{aligned} \mathcal{M}_k &= \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{expit}(c_0 + \mathbf{c}_1^\top \mathbf{x}) (2\pi)^{-d/2} \exp\left\{ \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x} \mathbf{x}^\top) \end{bmatrix}^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} d\mathbf{x} \\ &\times (2\pi)^{d/2} \exp(A(\boldsymbol{\eta})). \end{aligned}$$

Taking the inverse of the natural parameters as per equation (1.14) and implementing the change of variables $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u}$ leads to

$$\mathcal{M}_k = \int_{\mathbb{R}^d} (\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u})^{\otimes k} \text{expit}(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + (\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1)^\top \mathbf{u}) \phi_{\mathbf{I}}(\mathbf{u}) d\mathbf{u} Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + (d/2) \log(2\pi))$. Using equations (4.20) - (4.22), we can now obtain each of the moments using only univariate quadrature. Note that the integrals required are renowned for being numerically unstable. For stability, we rewrite the integrals in the forms of those presented by Kim & Wand (2017),³²

$$\mathcal{C}_b(k, r, q) = \int_{-\infty}^{\infty} x^k \exp(rx - qx^2 - b(x)) dx, \quad (4.24)$$

where $b(x) = \log(1 + \exp(x))$. Using simple algebraic manipulations it is easy to show

$$\begin{aligned} &\int_{-\infty}^{\infty} u^k \text{expit}(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1\|u) \phi(u) du \\ &= \int_{-\infty}^{\infty} ((r_2 + 2r_7x)(2r_7)^{-1/2})^k \exp(r_6x - r_7x^2 - b(x)) dx Z_0, \end{aligned}$$

where

$$r_1 = -2\mathbf{c}_1^\top \boldsymbol{\Sigma} \mathbf{c}_1, \quad r_2 = 2(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu})r_1^{-1}, \quad r_6 = 1 - r_2, \quad r_7 = -r_1^{-1}$$

and

$$Z_0 = \exp\left(\left(\frac{r_2}{2}\right)^2 r_1 + (1/2) \log(r_7/\pi)\right).$$

It follows,

$$\begin{aligned} Z_0^{-1} \int_{-\infty}^{\infty} \text{expit}(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1\| u) \phi(u) du &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} \int_{-\infty}^{\infty} u \text{expit}(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1\| u) \phi(u) du \\ &= (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7))(2r_7)^{-1/2}, \end{aligned}$$

and

$$\begin{aligned} Z_0^{-1} \int_{-\infty}^{\infty} u^2 \text{expit}(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1\| u) \phi(u) du \\ = 2(r_2 \mathcal{C}_b(1, r_6, r_7) + r_7 \mathcal{C}_b(2, r_6, r_7)) - (r_2^2 r_1 / 2) \mathcal{C}_b(0, r_6, r_7). \end{aligned}$$

It is then easy to find each of the k th moments

$$\begin{aligned} Z_0^{-1} Z_1^{-1} \mathcal{M}_0 &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_1 &= \boldsymbol{\mu} \mathcal{C}_b(0, r_6, r_7) + \boldsymbol{\Sigma} \mathbf{c}_1 (2r_7 \mathcal{C}_b(1, r_6, r_7) + r_2 \mathcal{C}_b(0, r_6, r_7)), \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_2 &= (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\Sigma}) \mathcal{C}_b(0, r_6, r_7) \\ &\quad + (\boldsymbol{\mu} \mathbf{c}_1^\top \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{c}_1 \boldsymbol{\mu}^\top) (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) \\ &\quad + \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma} (4r_7^2 \mathcal{C}_b(2, r_6, r_7) + 4r_7 r_2 \mathcal{C}_b(1, r_6, r_7) + (r_2^2 - 2r_7) \mathcal{C}_b(0, r_6, r_7)). \end{aligned}$$

Thus the optimal mean parameter for the projection is

$$\boldsymbol{\mu}^* = E(\mathbf{x}) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{c}_1 (r_2 + 2r_7 \mathcal{C}_{b,1:0}(r_6, r_7)),$$

and the optimal variance parameter is

$$\begin{aligned} \boldsymbol{\Sigma}^* &= E(\mathbf{x} \mathbf{x}^\top) - E(\mathbf{x}) E(\mathbf{x})^\top \\ &= \frac{\mathcal{M}_2}{\mathcal{M}_0} - \frac{\mathcal{M}_1}{\mathcal{M}_0} \left(\frac{\mathcal{M}_1}{\mathcal{M}_0} \right)^\top \\ &= \boldsymbol{\Sigma} + (2r_7)^2 \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma} (\mathcal{C}_{b,2:0}(r_6, r_7) - \mathcal{C}_{b,1:0}(r_6, r_7)^2 + r_1/2), \end{aligned}$$

where $\mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$ and $\mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$. By converting back to natural optimal and input parameters we arrive at Definition [17](#)

4.3.3 Proof of Definition 18

We wish to obtain the projection of an input function following the form of equation (4.12) onto the multivariate normal family. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. For the general case for all $\otimes k$. As before, we can obtain closed form solutions to the equivalent of the zeroth, first and second moments. Using the $\mathbf{x}^{\otimes k}$ notation as described in equation (1.2), we can write

$$\begin{aligned}\mathcal{M}_k &\equiv \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{expit}_{\text{MS}}(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp\left(\boldsymbol{\eta}_1^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2 \mathbf{x}\right) d\mathbf{x} \\ &\equiv \sum_{i=1}^8 p_i \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(s_i c_0 + s_i \mathbf{c}_1^\top \mathbf{x}) \exp\left(\boldsymbol{\eta}_1^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2 \mathbf{x}\right) d\mathbf{x},\end{aligned}$$

where $\text{expit}_{\text{MS}}(x)$ is given in equation (4.4). Using the matrix notation defined in equation (1.14) we can write

$$\begin{aligned}\mathcal{M}_k &= \sum_{i=1}^8 p_i \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Phi(s_i c_0 + s_i \mathbf{c}_1^\top \mathbf{x}) (2\pi)^{-d/2} \exp\left\{\left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta} - A(\boldsymbol{\eta})\right\} d\mathbf{x} \\ &\quad \times (2\pi)^{d/2} \exp(A(\boldsymbol{\eta})).\end{aligned}$$

Using the inverse of the natural parameters in equation (1.2) and implementing the change of variable $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u}$

$$\mathcal{M}_k = \sum_{i=1}^8 p_i \int_{\mathbb{R}^d} (\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u})^{\otimes k} \Phi\left(s_i(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu}) + (s_i \boldsymbol{\Sigma}^{1/2} \mathbf{c}_1)^\top \mathbf{u}\right) \phi_I(\mathbf{u}) d\mathbf{u} Z_1,$$

where $Z_1 = \exp\left(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi)\right)$. Then using Lemma 3 we can now obtain exact forms for each of the moments and thus calculate the optimal parameters for the projection. It is then easy to show each k th moment,

$$\begin{aligned}Z_1^{-1} \mathcal{M}_0 &= \sum_{i=1}^8 p_i \Phi(r_{2i}), \\ Z_1^{-1} \mathcal{M}_1 &= \boldsymbol{\mu} \sum_{i=1}^8 p_i \Phi(r_{2i}) + \boldsymbol{\Sigma} \mathbf{c}_1 \sum_{i=1}^8 (2p_i s_i r_{1i}^{-1}) \phi(r_{2i}), \\ Z_1^{-1} \mathcal{M}_2 &= (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\Sigma}) \sum_{i=1}^8 p_i \Phi(r_{2i}) + (\boldsymbol{\mu} \mathbf{c}_1^\top \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{c}_1 \boldsymbol{\mu}^\top) \sum_{i=1}^8 (2p_i s_i r_{2i}^{-1}) \phi(r_{2i}) \\ &\quad - \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma} \sum_{i=1}^8 (4p_i s_i^2 r_{2i} r_{1i}^{-2}) \phi(r_{2i}),\end{aligned}$$

where $r_{1i} = 2(s_i^2 \mathbf{c}_1^\top \boldsymbol{\Sigma} \mathbf{c}_1 + 1)^{1/2}$ and $r_{2i} = 2s_i(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu}) r_{1i}^{-1}$. The optimal mean and variance parameters for the projection follow and are respectively

$$\boldsymbol{\mu}^* = E(\mathbf{x}) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{c}_1 r_3$$

and

$$\boldsymbol{\Sigma}^* = E(\mathbf{x}\mathbf{x}^\top) - E(\mathbf{x})E(\mathbf{x})^\top = \frac{\mathcal{M}_2}{\mathcal{M}_0} - \frac{\mathcal{M}_1}{\mathcal{M}_0} \left(\frac{\mathcal{M}_1}{\mathcal{M}_0} \right)^\top = \boldsymbol{\Sigma} - 2r_4 \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma},$$

where

$$r_4 = \frac{1}{2}(\tilde{r}_3 + r_3^2), \quad r_3 = \frac{\sum_{i=1}^8 (2p_i s_i r_{1i}^{-1}) \phi(r_{2i})}{\sum_{i=1}^8 p_i \Phi(r_{2i})} \quad \text{and} \quad \tilde{r}_3 = \frac{\sum_{i=1}^8 (4p_i s_i^2 r_{2i} r_{1i}^{-2}) \phi(r_{2i})}{\sum_{i=1}^8 p_i \Phi(r_{2i})}.$$

By converting both the input and optimal parameters back to natural parameters we arrive at the desired result.

Chapter 5

Expectation propagation for one level count response mixed models

Thus far, the novel methodology presented in this thesis has concerned datasets with binary response variables. We now consider Poisson and negative binomial models for datasets with count response variables. As the projections required for EP cannot be solved via closed form solutions, we implement the quadrature approach used for the logistic models. The message passing framework implemented in the binary case remains mostly unchanged, with only minor alterations required. Additionally, care must be taken during optimisation of the shape parameter for the negative binomial model. We negate any repeated details of work presented in previous chapters and refer readers back where appropriate.

This chapter is broken into four main sections. Section [5.1](#) explains the simplest random intercepts only Poisson model, while Section [5.2](#) discusses its extension to a general model for any number of fixed and random effects. Section [5.3](#) explains a simple model for inference on the shape parameter of the negative binomial model, while Section [5.5](#) explains its extension to a general model for any number of fixed and random effects.

5.1 The simplest Poisson mixed model

We consider a random intercepts only model first for simplicity with a balanced dataset, where all m groups have the same n number of observations in them. Note for the random intercept only model the EP approach affords no benefit over directly computing the integrals via quadrature, acting as a test bed for the general case which follows. For observed values of

$$y_{ij}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

where $y_{ij} \in \mathbb{Z}_{\geq 0}$, the form of the model is

$$y_{ij}|u_i \stackrel{\text{iid}}{\sim} \text{Poisson}(\exp(u_i); y_{ij}), \quad u_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2),$$

where $\text{Poisson}(\lambda; x)$ is given by equation (1.11) and u_i is a scalar unobserved latent variable. We wish to find the maximiser of $\ell(\sigma^2)$ denoted by $\widehat{\sigma^2}$ and the best predictor of the random effects $\text{BP}(u_i)$.

The log-likelihood of the simplified model can be written as

$$\ell(\sigma^2) = \sum_{i=1}^m \ell_i(\sigma^2),$$

where

$$\ell_i(\sigma^2) \equiv \log \int_{-\infty}^{\infty} \prod_{j=1}^n \text{Poisson}(\exp(u_i); y_{ij}) (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/2\sigma^2) du_i, \quad (5.1)$$

and the maximum likelihood estimate of σ^2 is given by

$$\widehat{\sigma^2} = \underset{\sigma^2}{\text{argmax}} \ell(\sigma^2).$$

The best predictor of the random effect is given by

$$\text{BP}(u_i) = \frac{\int_{-\infty}^{\infty} u_i \prod_{j=1}^n \text{Poisson}(\exp(u_i); y_{ij}) \exp(-u_i^2/2\sigma^2) du_i}{\int_{-\infty}^{\infty} \prod_{j=1}^n \text{Poisson}(\exp(u_i); y_{ij}) \exp(-u_i^2/2\sigma^2) du_i}.$$

Using EP, we develop an approximation and sum each $\ell_i(\sigma^2)$ to obtain the full log-likelihood and best-predictors. We compare it to a traditional quadrature approach and show both methods provide reasonable and similar estimates.

The following subsections copy the structure of Chapter 4. In Subsection 5.1.1 we

provide details of the quadrature approach to estimating the likelihood surface, then explain our novel method using EP in Subsection 5.1.2. We compare the likelihood surface of both methods in Subsection 5.1.3, and outline computation of best predictors in Subsection 5.1.4. Point estimate and confidence interval calculation is conducted analogously to Section 2.4.1 and as such we do not repeat it.

5.1.1 Traditional quadrature likelihood approximation

Implementation of adaptive quadrature via Gauss-Kronrod follows the probit and exit cases explained in Sections 2.1 and 4.1.1 respectively, where we utilise the R function `integrate()` in the “stats” package (R Core Team, 2019⁵⁶). Denoting $y_i. \equiv \sum_{j=1}^n y_{ij}$, we suggest that for numerical stability each $\ell_i(\sigma^2)$ is calculated as

$$\frac{1}{2} \log(2\pi\sigma^2) + \ell_i(\sigma^2) = h(u_{0i}) + \log \int_{-\infty}^{\infty} \exp(h_i(u) - h_i(u_{0i})) du,$$

where

$$\begin{aligned} h_i(u) &\equiv y_i. u - n \exp(u) - \frac{u^2}{2\sigma^2}, \\ h'_i(u) &= y_i. - n \exp(u) - \frac{u}{\sigma^2} \end{aligned}$$

and u_{0i} is the root of h'_i , which we recommend finding by a bisection search, where the starting values are selected -1 and 1 to be for the lower and upper bounds respectively.

5.1.2 Expectation propagation likelihood approximation

We now consider an EP approach to the approximate likelihood (denoted by $\underline{\ell}(\sigma^2)$) as in Section 2.2. The EP approximation requires an unnormalised normal density function selected by minimising the KL-divergence criterion, to replace each

$$\text{Poisson}(\exp(u_i); y_{ij}), \quad 1 \leq j \leq n$$

in equation (5.1). For the case of Poisson GLMMs, the required projections onto an unnormalised univariate normal distribution can be written as

$$f_{\text{input}}(x) = \text{Poisson}(\exp(c_0 + c_1 x); c_2) \exp(\eta_1^{\text{input}} x + \eta_2^{\text{input}} x^2), \quad (5.2)$$

where $\eta_1^{\text{input}} \in \mathbb{R}$, $\eta_2^{\text{input}} < 0$, $c_0 = 0$, $c_1 = 1$, $c_2 = y_{ij}$ and $x = u_i$. Subsequently, the integrand is proportional to a product of univariate normal density functions. As such, we seek $\boldsymbol{\eta}^*$ such that

$$\begin{aligned} & \int_{-\infty}^{\infty} x^k \text{Poisson}(\exp(c_0 + c_1 x); c_2) \exp(\eta_1^{\text{input}} x + \eta_2^{\text{input}} x^2) dx \\ &= \int_{-\infty}^{\infty} x^k \exp \left\{ \left[\begin{array}{c} 1 \\ x \\ x^2 \end{array} \right]^\top \boldsymbol{\eta}^* \right\} dx. \end{aligned} \quad (5.3)$$

To obtain the required projection, we first obtain the optimal natural parameters η_1^* and η_2^* to project onto the normal family. Since the integral arising in the required projection does not have a closed form solution, we use univariate quadrature to obtain it. As before, caution must be exercised when calculating the integrals since they are prone to be numerically unstable. As such, each integral is calculated using the form of $C_b(p, q, r)$ as in Section 2.1 of Kim & Wand (2018)³² for numerical stability. By doing so we arrive at Result 20.

Result 20. Given f_{input} follows the form of equation (5.2), the projection onto the univariate normal family is given by

$$\text{proj}_N[f_{\text{input}}] = \exp \left(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*) \right) h(x),$$

where

$$\boldsymbol{\eta}_{-1}^* = k_{\text{Poisson}}(\boldsymbol{\eta}_{-1}^{\text{input}}; c_0, c_1, c_2), \quad \boldsymbol{\eta}_{-1}^{\text{input}} \equiv \begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix},$$

$k_{\text{poisson}} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1, c_2 \right)$ is defined in Definition 21 and $\mathbf{T}(x)$ and $h(x)$ follow from equation (1.5.2.1).

Definition 21. For primary scalar arguments $a_1 \in \mathbb{R}$ and $a_2 < 0$, and auxiliary scalar arguments $c_0, c_1 \in \mathbb{R}$ and $c_2 \in \mathbb{Z}_{\geq 0}$, the function $k_{\text{Poisson}} : H \rightarrow H$ is given by

$$k_{\text{Poisson}} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1, c_2 \right) \equiv \begin{bmatrix} r_5(a_1 + r_3 c_1) \\ r_5 a_2 \end{bmatrix}, \quad (5.4)$$

where

$$\begin{aligned} r_1 &= c_1^2 a_2^{-1}, & r_2 &= (c_1 a_2^{-1} a_1 - 2c_0) r_7, & r_3 &= r_2 + 2\mathcal{C}_{b,1:0}(r_6, r_7) r_7, \\ r_4 &= 2(\mathcal{C}_{b,1:0}(r_6, r_7))^2 - \mathcal{C}_{b,2:0}(r_6, r_7) - r_1 r_7^2, & r_5 &= (a_2 + r_4 c_1^2)^{-1} a_2, \\ r_6 &= 1 - r_2, & r_7 &= -(r_1)^{-1}, & \mathcal{C}_{b,1:0}(r_6, r_7) &= \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, \\ \mathcal{C}_{b,2:0}(r_6, r_7) &= \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)} \quad \text{and} \quad b(x) = \exp(x). \end{aligned}$$

A proof of Definition 21 is provided in Appendix 5.6.1. Using Result 20 we now obtain the normalising natural parameter η_0^* to find the projection onto unnormalised normal family.

5.1.2.1 Projection onto the unnormalised normal family

Recall the moment matching problem described by equation (2.6) and Result 4. Then the normalising factor can be shown to be

$$C_f = \int_{\mathbb{R}^d} f_{\text{input}}(x) dx = \mathcal{C}_b(0, r_6, r_7) Z_0 Z_1,$$

where $Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi))$, $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) - \log \Gamma(c_2 + 1))$ and r_6 and r_7 are given in Definition 21. By Result 4

$$\begin{aligned} \boldsymbol{\eta}_0^* &= \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^2 / \boldsymbol{\eta}_2^* \\ &\quad - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^2 / \boldsymbol{\eta}_2^{\text{input}} + \frac{1}{2} \log(\boldsymbol{\eta}_2^* / \boldsymbol{\eta}_2^{\text{input}}) - \log \Gamma(c_2 + 1). \end{aligned}$$

Thus to obtain η_0^* we introduce Definition 22.

Definition 22. For primary scalar arguments $a_1, b_1 \in \mathbb{R}$ and $a_2, b_2 < 0$ and auxiliary scalar arguments $c_0, c_1 \in \mathbb{R}$ and $c_2 \in \mathbb{Z}_{\geq 0}$ the function $c_{Poisson} : H \rightarrow \mathbb{R}$ is defined as:

$$c_{Poisson} \left(\begin{array}{c} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ ; c_0, c_1, c_2 \end{array} \right) \equiv \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) \\ + \frac{1}{4} a_1^2/a_2 - \frac{1}{4} b_1^2/b_2 + \frac{1}{2} \log(a_2/b_2) - \log \Gamma(c_2 + 1),$$

where r_1, r_2, r_6 and r_7 follow as defined in Definition [21](#).

In summary, the projection of the input function onto the unnormalised normal family is given in Result [21](#).

Result 21. For an unnormalised input function of the form of equation [\(5.2\)](#),

$$proj_{UN}[f_{input}] = \exp \left\{ \begin{array}{c} \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{bmatrix} \end{array} \right\},$$

where

$$\begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix} = k_{Poisson} \left(\begin{array}{c} \begin{bmatrix} \eta_1^{input} \\ \eta_2^{input} \end{bmatrix} \\ ; c_0, c_1, c_2 \end{array} \right)$$

and

$$\eta_0^* = c_{Poisson} \left(\begin{array}{c} \begin{bmatrix} \eta_1^{input} \\ \eta_2^{input} \end{bmatrix}, \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix} \\ ; c_0, c_1, c_2 \end{array} \right).$$

We now show how to implement the results shown in this section in a message passing framework.

5.1.2.2 Message passing formulation

As mentioned, only minor changes to Algorithm [3](#) are required to account for the Poisson case. Note the components of the sum of the log-likelihood function $\ell_i(\sigma^2)$ are

$$\ell_i(\sigma^2) = \log \int_{-\infty}^{\infty} \left(\prod_{j=1}^n p(y_{ij}|u_i) \right) p(u_i; \sigma^2) du_i, \quad (5.5)$$

where

$$p(y_{ij}|u_i) \equiv \text{Poisson}(\exp(u_i), y_{ij}) \quad \text{and} \quad p(u_i; \sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/(2\sigma^2))$$

are respectively the conditional density functions of each response given its random effect and the density function of the random effect.

Since the dependence structure of the product in equation (5.5) matches that of the probit model shown in Figure 2.1, we only need to update the messages depending on the conditional density functions of each response given its random effect. Following Section 2.2.2, this involves updating messages $m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i)$ as in equation (2.10) to

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}_{\text{UN}} [\text{Poisson}(\exp(c_0 + c_1 u_i); c_{2_{ij}}) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)]}{\exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)},$$

where $c_0 = 0$, $c_1 = 1$ and $c_{2_{ij}} = y_{ij}$. Obtaining the required projection follows from Result 21 analogous to equation (2.16), where the linear and quadratic coefficient updates in equation (2.17) are changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} \leftarrow k_{\text{Poisson}}(\boldsymbol{\eta}_{1:2}^\otimes; c_0, c_1, c_{2_{ij}}) - \boldsymbol{\eta}_{1:2}^\otimes \quad (5.6)$$

and where the constant coefficient update in equation (2.18) is changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 \leftarrow c_{\text{Poisson}}(\boldsymbol{\eta}_{1:2}^\otimes, (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_0, c_1, c_{2_{ij}}). \quad (5.7)$$

In summary, Algorithm 3 applies to the Poisson model, however equations (2.17) and (2.18) are replaced with equations (5.6) and (5.7), with constant term inputs set as $c_0 \leftarrow 0$, $c_1 \leftarrow 1$ and $c_{2_{ij}} \leftarrow y_{ij}$ for $1 \leq i \leq m$, $1 \leq j \leq n$. Additionally, the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$ follows equation (5.9).

5.1.2.3 Starting values for the Poisson case

We now derive starting values for the Poisson case using the same principles shown in the univariate probit and logistic case. Let

$$\log p(y_{ij}|u_i) = f(u_i, y_{ij}) - \log \Gamma(y_{ij} + 1),$$

where

$$\begin{aligned} f(x, y_{ij}) &= f(x) = y_{ij}x - \exp(x), \\ f'(x) &= y_{ij} - \exp(x), \end{aligned}$$

and

$$f''(x) = -\exp(x). \quad (5.8)$$

Let \hat{u}_i be an approximation of u_i . Then a Taylor series expansion of $f(u_i)$ evaluated at \hat{u}_i leads to

$$\begin{aligned} f(u_i) &= f(\hat{u}_i) + f'(\hat{u}_i)(u_i - \hat{u}_i) + \frac{1}{2}f''(\hat{u}_i)(u_i - \hat{u}_i)^2 + \dots \\ &= \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{u}_i) \\ f'(\hat{u}_i) \\ \frac{1}{2}f''(\hat{u}_i) \end{bmatrix}.$$

The quadratic approximation to $\log p(y_{ij}|u_i)$ based on Taylor expansion about \hat{u}_i is $\log \check{p}(y_{ij}|u_i)$ as per equation (2.22). By following the same logic of the previous cases it is easy to show

$$\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ f'(\hat{u}_i) - f''(\hat{u}_i)\hat{u}_i \\ \frac{1}{2}f''(\hat{u}_i) \end{bmatrix}, \quad (5.9)$$

where

$$\eta_0^{\text{start}} = f(\hat{u}_i) - f'(\hat{u}_i)\hat{u}_i + \frac{1}{2}f''(\hat{u}_i)(\hat{u}_i)^2.$$

The comments at the end of Section 2.2.3 apply for these starting values.

5.1.3 Evaluation of the estimates

Implementing the quadrature and EP approaches in the R computing environment, we now visually compare the accuracy of our likelihood approximation $\underline{\ell}'(\sigma^2)$ to the exact likelihood surface $\ell(\sigma^2)$. Figure 5.1 plots the estimates of the likelihood surface for both methods. The data generated had 50 groups with 5 responses per group (i.e. $m = 50$ and $n = 5$) and the true value of $\sigma^2 = 0.17$. Using the quadrature method as exact, the plot shows although there are some discrepancies on the tails of the likelihood surface, the approximate method follows the exact likelihood surface around the maximum well. Additionally, the true value of σ^2 matches well with the maximum of the exact likelihood.

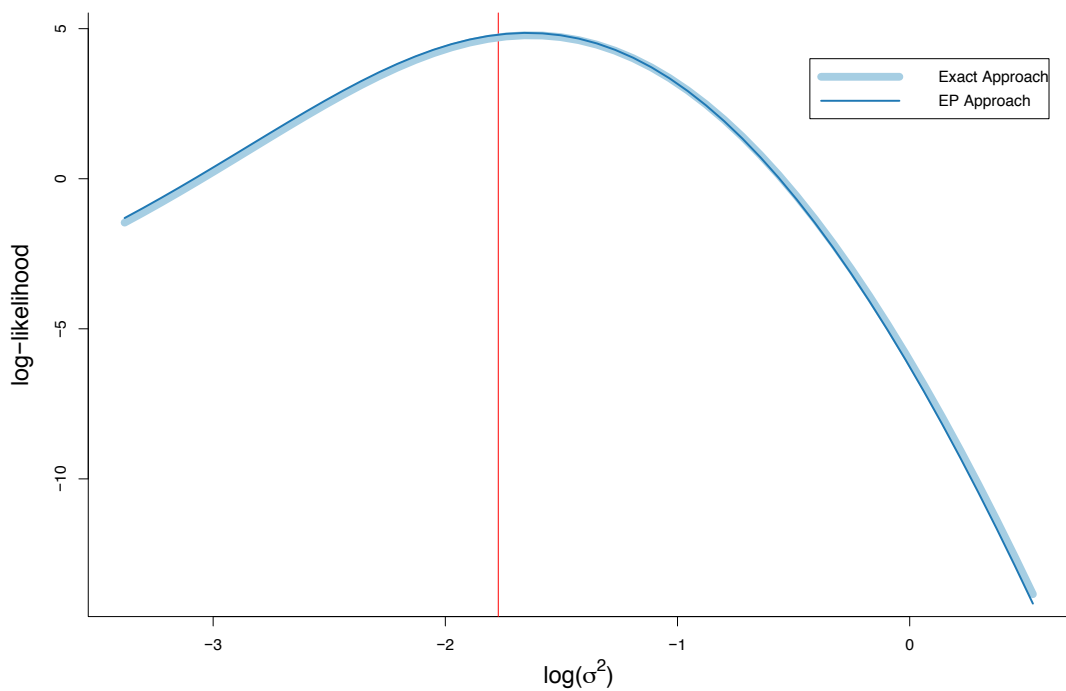


Figure 5.1: A comparison plot of the log-likelihood surface over the parameter σ^2 for Poisson models calculated exactly using univariate quadrature and approximated via EP. The true $\sigma^2 = 0.17$ is represented on the log scale by the red line. The EP approximation is shown by the dark blue line and the exact by the light blue line.

5.1.4 Best predictor

Best prediction for σ^2 via quadrature and EP is the same as presented in Section 2.5, although for the former we now redefine

$$J_{si}(\sigma^2) \equiv \int_{-\infty}^{\infty} u_i^s \prod_{j=1}^n \text{Poisson}(\exp(u_i); y_{ij}) (2\pi\sigma^2)^{-1/2} \exp(-u_i^2/(2\sigma^2)) du_i, \quad s \in \{0, 1, 2\}.$$

The calculations are otherwise analogous.

5.2 General Poisson mixed models

We now consider EP for the general Poisson model where any number of fixed and random effects can be specified. Consider

$$y_{ij} | \mathbf{u}_i \stackrel{\text{ind}}{\sim} \text{Poisson}\left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}); y_{ij}\right), \quad \mathbf{u}_i \stackrel{\text{ind}}{\sim} \mathbf{N}(\mathbf{0}_{d\mathbf{R}}, \boldsymbol{\Sigma}),$$

$$1 \leq i \leq m \quad \text{and} \quad 1 \leq j \leq n_i,$$

where the notation follows the general one level model presented in Section 1.8. The log-likelihood can be expressed as

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}),$$

where

$$\begin{aligned} \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \log \int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{Poisson}\left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}); y_{ij}\right) \right\} \\ &\quad \times |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i, \end{aligned}$$

and the best predictor of \mathbf{u}_i is

$$\text{BP}(\mathbf{u}_i) \equiv \frac{\int_{\mathbb{R}^{d\mathbf{R}}} \mathbf{u}_i \left\{ \prod_{j=1}^{n_i} \text{Poisson}\left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}); y_{ij}\right) \right\} \exp\left(-\frac{1}{2}\mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i}{\int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{Poisson}\left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}); y_{ij}\right) \right\} \exp\left(-\frac{1}{2}\mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) d\mathbf{u}_i}.$$

(5.10)

Implementation of EP follows as per the previous models, with minor alterations to the algorithm to account for the count response variable.

We first explain likelihood approximation using EP in Subsection 5.2.1, before results of a simulation study are presented in Subsection 5.2.2. Details regarding computation of point estimates and confidence intervals are analogous to those presented in Section 3.2, whilst the same calculations in Section 3.3 can be used for best prediction of \mathbf{u}_i . We do not repeat either section again and instead refer readers to the previous work.

5.2.1 Expectation propagation likelihood approximation

EP centres around finding the optimal natural parameters η_0 , η_1 and η_2 which minimise $\text{KL}(f_{\text{input}} \parallel f_{UN})$, where f_{UN} is defined by equation (3.4) and

$$f_{\text{input}}(\mathbf{x}) = \text{Poisson}\left(\exp(c_0 + \mathbf{c}_1^\top \mathbf{x}); c_2\right) \exp\left(\left(\boldsymbol{\eta}_1^{\text{input}}\right)^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x}\right), \quad (5.11)$$

where $\boldsymbol{\eta}_1^{\text{input}}$ is a $d \times 1$ vector, $\mathbf{H}_2^{\text{input}}$ is a $d \times d$ matrix, $c_0 = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}$, $c_2 = y_{ij}$, $\mathbf{c}_1 = \mathbf{x}_{ij}^{\mathbf{R}}$ and $x = \mathbf{u}_i$. As such, we seek an $\boldsymbol{\eta}^*$ to solve

$$\begin{aligned} & \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{Poisson}\left(\exp(c_0 + \mathbf{c}_1^\top \mathbf{x}); c_2\right) \exp\left\{\left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^{\text{input}}\right\} d\mathbf{x} \\ &= \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \exp\left\{\left[\begin{array}{c} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^*\right\} d\mathbf{x}, \end{aligned} \quad (5.12)$$

where $\mathbf{x}^{\otimes k}$ is as defined in equation (1.2). Thus to obtain the required projection, we first obtain the optimal natural parameters $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$ to project onto the multivariate normal family as is presented in Result 22. As the integrals required to do so are not available in closed form solutions, we use the same properties as the exploit for Poisson models to express the multivariate integrals as univariate integrals. These integrals are expressed as in Section 2.1 of Kim & Wand (2018)³² as denoted by $\mathcal{C}_b(p, q, r)$ and solved using quadrature in a more efficient manner.

Result 22. Given f_{input} follows the form of equation (5.11), the projection onto the multivariate normal family is given by

$$proj_{\mathcal{N}}[f_{input}] = \exp\left(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*)\right)h(\mathbf{x}),$$

where

$$\boldsymbol{\eta}_{-1}^* \equiv K_{Poisson}(\boldsymbol{\eta}_{-1}^{input}; c_0, \mathbf{c}_1, c_2),$$

$$\boldsymbol{\eta}_{-1}^{input} \equiv \begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix},$$

$K_{Poisson}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2\right)$ is defined in Definition 23 and $\mathbf{T}(\mathbf{x})$ and $h(\mathbf{x})$ follow from Section 1.5.2.2.

Definition 23. For primary arguments \mathbf{a}_1 ($d \times 1$) and \mathbf{a}_2 ($\frac{1}{2}d(d+1) \times 1$) such that $vec^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{a}_2)$ is symmetric and positive definite, and auxiliary arguments $c_0 \in \mathbb{R}$, $c_2 \in \mathbb{Z}^+$ and \mathbf{c}_1 ($d \times 1$), the function $K_{Poisson} : H \rightarrow H$ is given by

$$K_{Poisson}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2\right) \equiv \begin{bmatrix} \mathbf{R}_5^\top (\mathbf{a}_1 + r_3 \mathbf{c}_1) \\ \mathbf{D}_d^\top vec(\mathbf{R}_5^\top \mathbf{A}_2) \end{bmatrix}, \quad (5.13)$$

where

$$\mathbf{A}_2 \equiv vec^{-1}((\mathbf{D}_d^+)^\top \mathbf{a}_2), \quad r_1 = \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{c}_1, \quad r_2 = (\mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1 - 2c_0)r_7,$$

$$r_3 = r_2 + 2 \mathcal{C}_{b,1:0}(r_6, r_7)r_7, \quad r_4 = 2(\mathcal{C}_{b,1:0}(r_6, r_7)^2 - \mathcal{C}_{b,2:0}(r_6, r_7) - r_1)r_7^2,$$

$$\mathbf{R}_5 = (\mathbf{A}_2 + r_4 \mathbf{c}_1 \mathbf{c}_1^\top)^{-1} \mathbf{A}_2, \quad r_6 = 1 - r_2, \quad r_7 = -(r_1)^{-1},$$

$$\mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, \quad \mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, \quad b(x) = \exp(x)$$

and $\mathcal{C}_b(k, r, q)$ is as per equation (5.42).

A proof of Definition 23 is given in Appendix 5.6.2. Using Result 22 we now obtain the normalising natural parameter $\boldsymbol{\eta}_0^*$ to find the projection onto unnormalised normal family.

5.2.1.1 Projection onto the unnormalised multivariate normal family

Recall the moment matching problem from equation (5.12) and Results 12 and 13. Then, we require

$$C_{f_{\text{input}}} = \int_{\mathbb{R}^d} f_{\text{input}}(\mathbf{x}) d\mathbf{x} = C_b(0, r_6, r_7) Z_0 Z_1,$$

where $Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi))$, $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) - \log \Gamma(c_2 + 1))$ and r_6 and r_7 are given in Definition 23. Analogous to previous arguments, we then get

$$\begin{aligned} \boldsymbol{\eta}_0^* &= \log C_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^\top (\mathbf{H}_2^*)^{-1} \boldsymbol{\eta}_1^* \\ &\quad - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^\top (\mathbf{H}_2^{\text{input}})^{-1} \boldsymbol{\eta}_1^{\text{input}} + \frac{1}{2} \log(|\mathbf{H}_2^*|/|\mathbf{H}_2^{\text{input}}|) - \log \Gamma(c_2 + 1). \end{aligned}$$

To find the normalising constant of the input function, we introduce Definition 24.

Definition 24. Consider first, primary arguments \mathbf{a}_1 and \mathbf{b}_1 and auxiliary argument \mathbf{c}_1 where all three are $d \times 1$. Next consider arguments \mathbf{a}_2 and \mathbf{b}_2 which are all $(\frac{1}{2}d(d+1) \times 1)$ such that both $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{a}_2)$ and $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{b}_2)$ are symmetric and positive definite. Finally note auxiliary scalar argument $c_0, \in \mathbb{R}$, $c_2, \in \mathbb{Z}^+$. Then the function $C_{\text{Poisson}} : H \times H \rightarrow \mathbb{R}$ is given by

$$\begin{aligned} C_{\text{Poisson}} \left(\begin{pmatrix} \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2 \end{pmatrix} \right) &\equiv \log C_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) \\ &\quad + \frac{1}{4} \mathbf{b}_1^\top \mathbf{B}_2^{-1} \mathbf{b}_1 - \frac{1}{4} \mathbf{a}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1 + \frac{1}{2} \log(|\mathbf{B}_2|/|\mathbf{A}_2|) - \log \Gamma(c_2 + 1), \end{aligned}$$

where $\mathbf{A}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{a}_2)$, $\mathbf{B}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{b}_2)$, r_1, r_2, r_6 and r_7 follow from Definition 23.

In summary, the projection onto the unnormalised multivariate normal family is given by Result 23.

Result 23. For an unnormalised input function of the form of equation (5.11),

$$proj_{\mathcal{UN}} [f_{input}] (\mathbf{x}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ vec(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} \right\},$$

where

$$\begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} = K_{Poisson} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}; c_0, \mathbf{c}_1, c_2 \right)$$

and

$$\eta_0^* = C_{Poisson} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1, c_2 \right).$$

5.2.1.2 Message passing formulation

As mentioned throughout this section, only minor changes to Algorithm 6 are required to account for the Poisson case. Each component of the sum of the log-likelihood function $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ can be written as

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \log \int_{\mathbb{R}^d} \left(\prod_{j=1}^{n_i} p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) d\mathbf{u}_i, \quad (5.14)$$

where

$$p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \equiv \text{Poisson} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), y_{ij} \right)$$

and

$$p(\mathbf{u}_i; \boldsymbol{\Sigma}) \equiv |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right)$$

are respectively the conditional density functions of each response given its random effect and the density function of that random effect.

As the dependence structure of the product in equation (5.14) is the same as the probit model shown in Figure 3.1, we only need to update the messages depending on the conditional density functions of each response given its random effect. This involves

updating messages $m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}(\mathbf{u}_i)$ as in equation (3.12)

$$m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}_{\mathcal{U}\mathcal{N}} \left[\text{Poisson} \left(\exp(c_{0_{ij}} + \mathbf{c}_{1_{ij}}^\top \mathbf{u}_i), c_{2_{ij}} \right) \exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\} \right]}{\exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\}},$$

where we set the constant terms:

$$c_{0_{ij}} \leftarrow \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}; \quad \mathbf{c}_{1_{ij}} \leftarrow \mathbf{x}_{ij}^{\mathbf{R}}; \quad c_{2_{ij}} \leftarrow y_{ij}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i.$$

Following Section 3.1.2, this task reduces down to adjusting the calculation of the optimal natural parameters $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ in equation (3.18) for the Poisson model, where the linear and quadratic coefficient updates in equation (3.19) given by K_{probit} are changed to

$$\left(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i} \right)_{1:2} \leftarrow K_{\text{Poisson}} \left(\boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ij}}, \mathbf{c}_{1_{ij}}, c_{2_{ij}} \right) - \boldsymbol{\eta}_{1:2}^\otimes \quad (5.15)$$

as per Definition 23, and where the constant coefficient update in equation (3.20) given by C_{probit} is changed to

$$\left(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i} \right)_0 \leftarrow C_{\text{Poisson}} \left(\boldsymbol{\eta}_{1:2}^\otimes, \left(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i} \right)_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ij}}, \mathbf{c}_{1_{ij}}, c_{2_{ij}} \right)$$

as per Definition 24.

Barring the aforementioned changes and the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ with the values discussed in Section 5.2.1.3, the full algorithm for the approximation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ follows Algorithm 6.

5.2.1.3 Starting values for the Poisson case

We now derive starting values for the Poisson case using the same principles shown in the probit and expit case. Let

$$\log p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) = f(a_{ij}; y_{ij}) - \log \Gamma(y_{ij} + 1),$$

where $a_{ij} = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}$ and $f(x)$ is defined as in equation (5.8). Then

$$(a_{ij} - \hat{a}_{ij}) = (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}}$$

where $\hat{a}_{ij} = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \hat{\mathbf{u}}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}$ and $\hat{\mathbf{u}}$ is an approximation of \mathbf{u}_i . A Taylor series expansion of $f(a_{ij})$ evaluated at \hat{a}_{ij} leads to

$$\begin{aligned} f(a_{ij}) &= f(\hat{a}_{ij}) + f'(\hat{a}_{ij})(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}} + \frac{1}{2} f''(\hat{a}_{ij})((\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}})^2 + \dots \\ &= \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top) \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{a}_{ij}) \\ f'(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} \\ \frac{1}{2} f''(\hat{a}_{ij}) \mathbf{D}_{d^{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix}.$$

Quadratic approximation to $\log p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})$ based on Taylor expansion about $\hat{\mathbf{u}}_i$ is $\log \check{p}(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})$ as per equation (3.22). By following the same logic of the previous cases it is easy to show

$$\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ f'(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} - f''(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i \\ \frac{1}{2} f''(\hat{a}_{ij}) \mathbf{D}_{d^{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix},$$

where

$$\eta_0^{\text{start}} = f(\hat{a}_{ij}) - f'(\hat{a}_{ij}) (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i + \frac{1}{2} f''(\hat{a}_{ij}) \hat{\mathbf{u}}_i^\top \mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i.$$

5.2.2 Simulation study

We now provide results of a simulation study for a random intercept model with 1000 replicates. Datasets were simulated with true parameter values

$$\boldsymbol{\beta}_{\text{true}} = [0.38, 0.93]^\top \quad \text{and} \quad \sigma_{\text{true}}^2 = -0.53, \quad (5.16)$$

generated with 50 groups and 5 measurements per group (i.e. $m = 50$ and $n = 5$). The $\mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\mathbf{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}]^\top \quad \text{and} \quad \mathbf{x}_{ij}^{\mathbf{R}} = 1,$$

where $x_{1_{ij}}$ was generated independently from a uniform distribution on the unit interval. The tolerance of error values were set to 10^{-5} for the EP scheme and the maximum number of iterations for optimisation was set to 1000.

We compare our EP methodology to Laplace approximation and adaptive Gauss-Hermite quadrature using 100 quadrature points. Both alternative methods were implemented via the R function `glmer()` from the R package “lme4” (Bates, et al., 2018^[5]). Since in the Poisson case quadrature is required for each projection of EP, EP is considerably slower than Laplace approximations and Gauss-Hermite quadrature. Additionally we note that to run the simulations within a reasonable time frame, they were separated and run across multiple high performance computers with of varying specification.

The resulting estimates and 95% confidence intervals for each interpretable model parameter of the study are presented in Figure 4.3, where the upper-right hand corner of each panel shows the empirical coverage values based on all 1000 replicates. Only 20 randomly chosen replicates from each method are shown in the panels for ease of viewing, where Laplace approximation and EP are shown by black lines, super imposed on grey lines showing the quadrature approach. Across the fixed effects parameters the empirical coverage of EP is within 0.2% of quadrature, where as Laplace approximations are marginally lower, particularly for the fixed slope which had coverage of 93.6%. The EP empirical coverage for the variance parameter was 1.6% higher than the expected 95%, where Laplace and quadrature are much closer coverage to 95%, with coverages of 94.9% and 95.1% respectively. Although we are aware that a comparison of time is obscured by factors such as language of implementation and computer performance, the average time for the EP routine over 100 replications was 1166.0054 seconds, whilst Laplace approximations and adaptive Gauss-Hermite took an average of 0.3108 and 0.3133 seconds respectively.

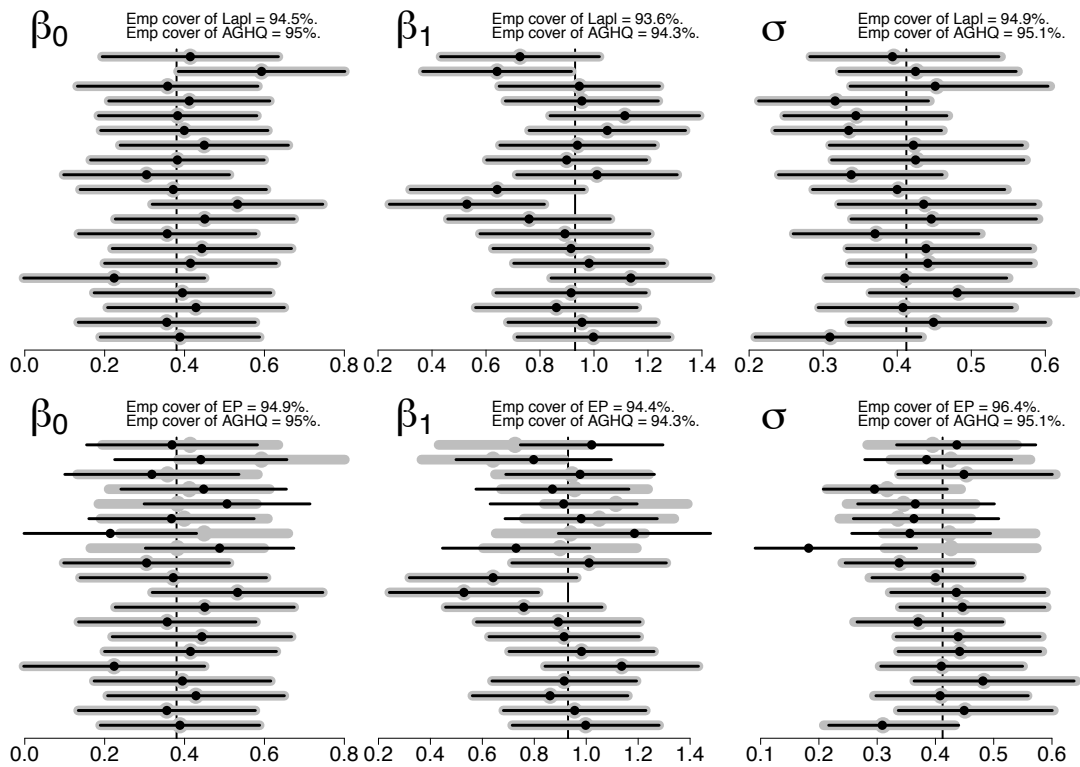


Figure 5.2: Summary comparison between confidence interval coverage for the univariate model with true parameter values from equation (5.16) for models fit with adaptive Gauss-Hermite quadrature and Laplace approximation. The horizontal lines are the confidence intervals for 20 randomly chosen replications of the simulation study, the solid circular points indicate the corresponding point estimates, the vertical lines indicate true parameter values, the grey lines correspond to adaptive Gauss-Hermite quadrature, the black lines on the top row correspond to Laplace approximation approach and the black lines on the bottom row correspond to EP approximation approach. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

5.3 The simplest negative binomial models

The negative binomial model follows naturally from the Poisson model. It provides the added benefit of a shape parameter κ , which facilitates handling over-dispersed count data. Specifically, as κ approaches zero the variance of the data becomes infinite. As κ increases to infinity the negative binomial distribution begins to resemble the Poisson distribution. Large values of κ can inturn lead to numerical issues, making the negative binomial model a particularly tricky one to handle. Analysts should determine whether to use Poisson or negative binomial links based on the nature of the data.

Since the negative binomial model brings the additional challenge of estimating the shape parameter κ , we begin developing our model on the case where we consider only the parameter κ . For this model we fix the variance between groups σ_{Fixed}^2 and aim to estimate κ . Additionally, we assume a balanced dataset with m -groups and n -observations per group. For observed values of

$$y_{ij}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n,$$

where $y_{ij} \in \mathbb{Z}_{\geq 0}$, the negative binomial model is

$$y_{ij}|u_i; \kappa \stackrel{\text{iid}}{\sim} \text{NB}(\exp(u_i), \kappa; y_{ij}), \quad u_i \stackrel{\text{iid}}{\sim} N(0, \sigma_{\text{Fixed}}^2),$$

where $\sigma_{\text{Fixed}}^2 \in \mathbb{R}_{>0}$ is a fixed number, $\kappa \in \mathbb{R}_{>0}$, $\text{NB}(y, \mu, \kappa)$ is given by equation (1.9) and u_i is a scalar unobserved latent variable. We wish to find the maximiser of $\ell(\kappa)$ denoted by $\hat{\kappa}$. Note that we do not find best predictors of u_i as they are assumed to be known in this setting. The likelihood of model can be expressed as

$$\ell(\kappa) = \sum_{i=1}^m \ell_i(\kappa),$$

where

$$\ell_i(\kappa) \equiv \log \int_{-\infty}^{\infty} \prod_{j=1}^n \text{NB}(\exp(u_i), \kappa; y_{ij}) (2\pi\sigma_{\text{Fixed}}^2)^{-1/2} \exp(-u_i^2/2\sigma_{\text{Fixed}}^2) du_i, \quad (5.17)$$

and the maximum likelihood estimate of κ is given by

$$\hat{\kappa} = \underset{\kappa}{\text{argmax}} \ell(\kappa).$$

Given the integral arising in the calculation of the log-likelihood surface of $\ell(\kappa)$ does not have a closed form solution, we develop an EP scheme to approximate it and compare the result to a quadrature approach.

Subsection [5.3.1](#) provides details of the quadrature approach to estimating the likelihood surface. We then explain our novel method using EP in Subsection [5.3.2](#), before comparing the likelihood surface of both methods in Subsection [5.3.3](#) and explaining point estimate and confidence interval calculation in Section [5.3.4](#).

5.3.1 Traditional quadrature likelihood approximation

Implementation of the adaptive quadrature via the Gauss-Kronrod method follows the same approach as in the Poisson case, where we utilise the R function `integrate()` in the “stats” package (R Core Team, 2019^[56]). We suggest that for numerical stability each $\ell_i(\kappa)$ arising in equation [\(5.17\)](#) is calculated as

$$\ell_i(\kappa) = h(u_{0i}) + \log \int_{-\infty}^{\infty} \exp(h_i(u) - h_i(u_{0i})) du,$$

where

$$\begin{aligned} h_i(u) &\equiv y_i u - \frac{u^2}{2\sigma_{\text{Fixed}}^2} - (n\kappa + y_i) \log(\kappa + \exp(u)), \\ h'_i(u) &\equiv y_i - \frac{u}{\sigma_{\text{Fixed}}^2} - \frac{(n\kappa + y_i) \exp(u)}{\log(\kappa + \exp(u))} \end{aligned}$$

and u_{0i} is the root of h'_i . We recommend finding u_{0i} using a bisection search, where the starting values are selected -1 and 1 to be for the lower and upper bounds respectively.

5.3.2 Expectation propagation likelihood approximation

We now consider an EP approach to the approximate likelihood denoted by $\underline{\ell}(\kappa)$. In the negative binomial case, the EP approximation requires an unnormalised normal density function to replace each

$$\text{NB}(\exp(u_i), \kappa; y_{ij}), \quad 1 \leq j \leq n$$

in equation [\(5.17\)](#), such that the KL-divergence criterion is minimised. For the case of negative binomial GLMMs, the required projections onto an unnormalised univariate

normal distribution are

$$f_{\text{input}}(x) = \text{NB}(\exp(c_0 + c_1x), \kappa; c_2) \exp(\eta_1^{\text{input}}x + \eta_2^{\text{input}}x^2), \quad (5.18)$$

where $\eta_1^{\text{input}} \in \mathbb{R}$, $\eta_2^{\text{input}} < 0$, $c_0 = 0$, $c_1 = 1$, $c_2 = y_{ij}$, $\kappa \in \mathbb{R}_{\geq 0}$ and $x = u_i$. Subsequently, the integrand is proportional to a product of univariate normal density functions. As such, we seek $\boldsymbol{\eta}^*$ such that

$$\begin{aligned} \int_{-\infty}^{\infty} x^k \text{NB}(\exp(c_0 + c_1x), \kappa; c_2) \exp(\eta_1^{\text{input}}x + \eta_2^{\text{input}}x^2) dx \\ = \int_{-\infty}^{\infty} x^k \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta}^* \right\} dx, \end{aligned} \quad (5.19)$$

Thus to obtain the required projection, we first obtain the optimal natural parameters η_1^* and η_2^* to project onto the normal family. Since the integral arising in the required projection does not have a closed form solution, we use univariate quadrature to obtain it. As before, it is important to use caution when calculating the integrals, since they are prone to be numerically unstable. For this reason we express each integral as in Section 2.1 of Kim & Wand (2018)³² as denoted by $\mathcal{C}_b(p, q, r)$. Using algebra similar to that required for the Poisson model, we arrive at Result [24](#).

Result 24. *Given f_{input} follows the form of equation [\(5.18\)](#), the projection onto the univariate normal family is given by*

$$\text{proj}_N[f_{\text{input}}] = \exp\left(\mathbf{T}(x)^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*)\right)h(x),$$

where

$$\boldsymbol{\eta}_{-1}^* = k_{\text{NB}}(\boldsymbol{\eta}_{-1}^{\text{input}}; c_0, c_1, c_2, \kappa), \quad \boldsymbol{\eta}_{-1}^{\text{input}} \equiv \begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix},$$

$k_{\text{NB}}\left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1, c_2, \kappa\right)$ is defined in Definition [25](#) and $\mathbf{T}(x)$ and $h(x)$ follow from Section [1.5.2.1](#)

Definition 25. For primary scalar arguments $a_1 \in \mathbb{R}$ and $a_2 \in \mathbb{R}_{\geq 0}$ and auxiliary scalar arguments $c_0, c_1 \in \mathbb{R}$, $c_2 \in \mathbb{Z}_{\geq 0}$ and $\kappa \in \mathbb{R}_{\geq 0}$ the function $k_{NB} : H \rightarrow H$ is given by

$$k_{NB} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}; c_0, c_1, c_2, \kappa \right) \equiv \begin{bmatrix} r_5(a_1 + r_3 c_1) \\ r_5 a_2 \end{bmatrix},$$

where

$$\begin{aligned} r_1 &= c_1^2 a_2^{-1}, & r_2 &= (c_1 a_2^{-1} a_1 - 2c_0) r_7, & r_3 &= r_2 + 2 \mathcal{C}_{b,1:0}(r_6, r_7) r_7, \\ r_4 &= 2(\mathcal{C}_{b,1:0}(r_6, r_7)^2 - \mathcal{C}_{b,2:0}(r_6, r_7) - r_1) r_7^2, & r_5 &= (a_2 + r_4 c_1^2)^{-1} a_2, \\ r_6 &= 1 - r_2, & r_7 &= -r_1^{-1}, & \mathcal{C}_{b,1:0}(r_6, r_7) &= \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, \\ \mathcal{C}_{b,2:0}(r_6, r_7) &= \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}, & b(x) &= (c_2 + \kappa) \log(\exp(x) + \kappa). \end{aligned}$$

and $\mathcal{C}_b(k, r, q)$ follows from equation (5.42).

A proof of Definition 25 is given in Appendix 5.6.3. Using Result 24 we now obtain the normalising natural parameter η_0^* to find the projection onto unnormalised normal family.

5.3.2.1 Projection onto the unnormalised normal family

Recall the moment matching problem described by equation (5.19) and Result 4. Then the normalising factor can be shown to be

$$C_f = \int_{\mathbb{R}^d} f_{\text{input}}(x) dx = \mathcal{C}_b(0, r_6, r_7) Z_0 Z_1,$$

where

$$Z_1 = \exp \left(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) + \log(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa) \right),$$

$Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi))$, and r_6 and r_7 are given in Definition 25. By Result 4

$$\begin{aligned} \eta_0^* &= \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} (\eta_1^*)^2 / \eta_2^* - \frac{1}{4} (\eta_1^{\text{input}})^2 / \eta_2^{\text{input}} \\ &\quad + \frac{1}{2} \log(\eta_2^* / \eta_2^{\text{input}}) + \log \Gamma(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa). \end{aligned}$$

Definition 26 follows.

Definition 26. For primary scalar arguments $a_1 \in \mathbb{R}$ and $a_2 \in \mathbb{R}_{\geq 0}$ and auxiliary scalar arguments $c_0, c_1 \in \mathbb{R}$, $c_2 \in \mathbb{Z}_{\geq 0}$ and $\kappa \in \mathbb{R}_{\geq 0}$ the function $c_{NB} : H \rightarrow \mathbb{R}$ is defined as:

$$\begin{aligned} c_{NB} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}; c_0, c_1, c_2, \kappa \right) &\equiv \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} a_1^2 / a_2 \\ &\quad - \frac{1}{4} b_1^2 / b_2 + \frac{1}{2} \log(a_2/b_2) + \log \Gamma(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa), \end{aligned}$$

where r_1, r_2, r_6 and r_7 follows from Definition 25.

Following the work of this section, the projection of the input function onto the unnormalised normal family is given in Result 25.

Result 25. For an unnormalised input function of the form of equation (2.5),

$$\text{proj}_{UN} [f_{\text{input}}] = \exp \left\{ \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \eta_1^* \\ \eta_2^* \end{bmatrix} \right\},$$

where

$$\begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix} = k_{NB} \left(\begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}; c_0, c_1, c_2, \kappa \right)$$

and

$$\eta_0^* = c_{NB} \left(\begin{bmatrix} \eta_1^{\text{input}} \\ \eta_2^{\text{input}} \end{bmatrix}, \begin{bmatrix} \eta_1^* \\ \eta_2^* \end{bmatrix}; c_0, c_1, c_2, \kappa \right).$$

We now show how to implement the results shown in this section in a message passing framework.

5.3.2.2 Message passing formulation

Only minor changes to Algorithm 3 are required to account for the negative binomial case. The components of the sum of the log-likelihood function $\ell_i(\kappa)$ are

$$\ell_i(\kappa) = \log \int_{-\infty}^{\infty} \left(\prod_{j=1}^n p(y_{ij}|u_i, \kappa) \right) p(u_i; \sigma_{\text{Fixed}}^2) du_i, \quad (5.20)$$

where

$$p(y_{ij}|u_i, \kappa) \equiv \text{NB}(\exp(u_i), \kappa; y_{ij})$$

and

$$p(u_i; \sigma_{\text{Fixed}}^2) \equiv (2\pi\sigma_{\text{Fixed}}^2)^{-1/2} \exp(-u_i^2/(2\sigma_{\text{Fixed}}^2))$$

are respectively the conditional density functions of each response given its random effect and the density function of that random effect.

Since the dependence structure of the product in equation (5.5) matches that of the probit model shown in Figure 2.1, we only need to update the messages depending on the conditional density functions of each response given its random effect. Following Section 2.2.2, this involves updating messages $m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i)$ as in equation (2.10)

$$m_{p(y_{ij}|u_i) \rightarrow u_i}(u_i) \leftarrow \frac{\text{proj}_{\text{UN}} [\text{NB}(\exp(c_0 + c_1 u_i), \kappa; c_{2_{ij}}) \exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)]}{\exp(\eta_1^\otimes u_i + \eta_2^\otimes u_i^2)},$$

where $c_0 = 0$, $c_1 = 1$, $c_{2_{ij}} = y_{ij}$ and $\kappa \in \mathbb{R}_{\geq 0}$. Utilising Result 21 leads to equation (2.16), where the linear and quadratic coefficient updates in equation (2.17) are changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} \leftarrow k_{\text{NB}}(\boldsymbol{\eta}_{1:2}^\otimes; c_0, c_1, c_{2_{ij}}, \kappa) - \boldsymbol{\eta}_{1:2}^\otimes \quad (5.21)$$

and where the constant coefficient update in equation (2.18) is changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_0 \leftarrow c_{\text{NB}}(\boldsymbol{\eta}_{1:2}^\otimes, (\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i})_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_0, c_1, c_{2_{ij}}, \kappa). \quad (5.22)$$

In summary, Algorithm 3 applies to the Poisson model, however equations (2.16) and (2.18) are replaced with equations (5.21) and (5.22), with constant term inputs set as $c_0 \leftarrow 0$, $c_1 \leftarrow 1$ and $c_{2_{ij}} \leftarrow y_{ij}$ for $1 \leq i \leq m$, $1 \leq j \leq n$. Additionally, the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}$ follows equation (5.24).

5.3.2.3 Starting values for the negative binomial case

We now derive starting values for the negative binomial case using the same principles shown in the previous cases. Let

$$\log p(y_{ij}|u_i, \kappa) = f(u_i, y_{ij}) + \log \Gamma(y_{ij} + \kappa) - \log \Gamma(y_{ij} + 1) - \log \Gamma(\kappa) + \kappa \log \kappa,$$

where

$$\begin{aligned} f(y_{ij}, x, \kappa) &= f(x) = y_{ij}x - (y + \kappa) \log(\exp(x) + \kappa), \\ f'(x) &= y_{ij} - \frac{(y + \kappa) \exp(x)}{\exp(x) + \kappa}, \end{aligned}$$

and

$$f''(x) = -\frac{\kappa(y + \kappa) \exp(x)}{(\exp(x) + \kappa)^2}. \quad (5.23)$$

Let \hat{u}_i be an approximation of u_i . Then a Taylor series expansion of $f(u_i)$ evaluated at \hat{u}_i leads to

$$\begin{aligned} f(u_i) &= f(\hat{u}_i) + f'(\hat{u}_i)(u_i - \hat{u}_i) + \frac{1}{2}f''(\hat{u}_i)(u_i - \hat{u}_i)^2 + \dots \\ &= \begin{bmatrix} 1 \\ u_i - \hat{u}_i \\ (u_i - \hat{u}_i)^2 \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{u}_i) \\ f'(\hat{u}_i) \\ \frac{1}{2}f''(\hat{u}_i) \end{bmatrix}.$$

By following the same logic of the probit case it is easy to show

$$\boldsymbol{\eta}_{p(y_{ij}|u_i) \rightarrow u_i}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ f'(\hat{u}_i) - f''(\hat{u}_i)\hat{u}_i \\ \frac{1}{2}f''(\hat{u}_i) \end{bmatrix}, \quad (5.24)$$

where

$$\eta_0^{\text{start}} = f(\hat{u}_i) - f'(\hat{u}_i)\hat{u}_i + \frac{1}{2}f''(\hat{u}_i)(\hat{u}_i)^2.$$

The comments at the end of Section [2.2.3](#) are also applicable for these starting values.

5.3.3 Evaluation of the estimates

Implementing the quadrature and EP approaches in the R computing environment, we now visually compare the accuracy of our likelihood approximation $\underline{\ell}'(\kappa)$ to the exact likelihood surface $\ell(\kappa)$. Figure [5.3](#) plots the estimates of the likelihood surface for both methods. The data generated had 100 groups with 5 responses per group and the true value of $\kappa = 15$ and $\sigma^2 = 0.25$. Using the quadrature method as exact, the plot shows that although there are some discrepancies on the tails of the likelihood surface, the approximate method follows the exact likelihood surface around the true maximum well. Additionally, the true value of κ matches well with the maximum of the exact likelihood.

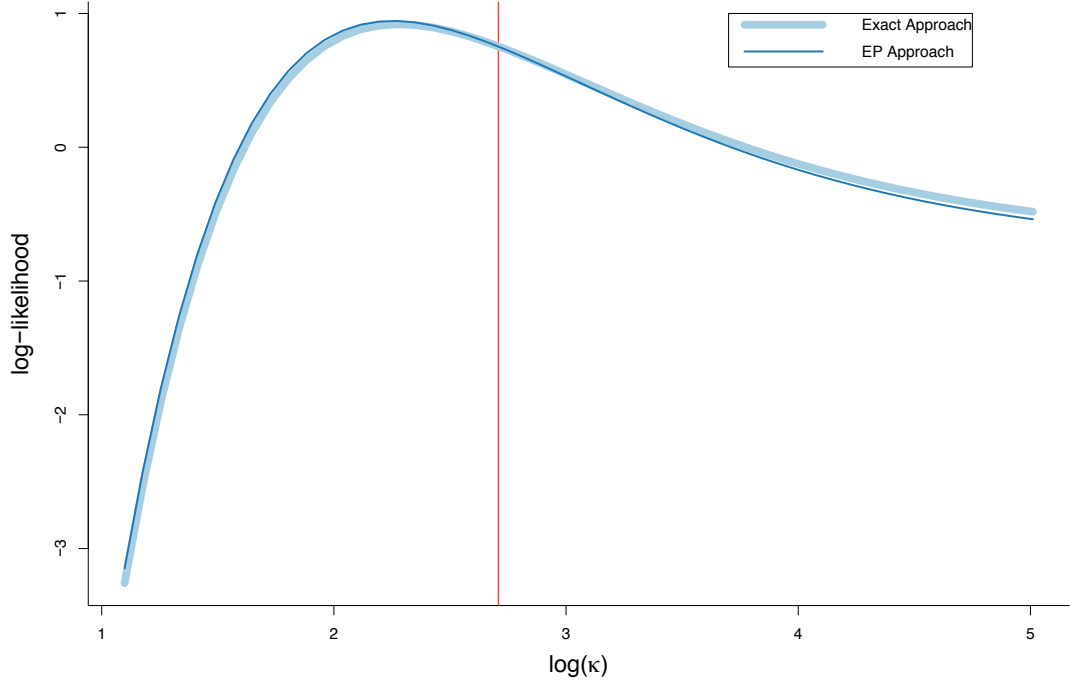


Figure 5.3: A comparison plot of the log-likelihood surface over the parameter κ for negative binomial models calculated exactly using quadrature and approximated via EP. The true $\kappa = 15$ is represented on the log scale by the red line. The EP approximation is shown by the dark blue line and the exact by the light blue line.

5.3.4 Computation of point estimates and confidence intervals

The maximum likelihood estimator for κ via quadrature and EP are respectively given by $\hat{\kappa} = \underset{\kappa}{\operatorname{argmax}} \ell(\kappa^2)$ and $\hat{\kappa} = \underset{\kappa}{\operatorname{argmax}} \underline{\ell}(\kappa)$. As before, to find their stationary points we require the first derivative of the likelihood functions denoted by $\ell'(\kappa)$ and $\underline{\ell}'(\kappa)$. Calculation of the second derivative, denoted by $\ell''(\kappa)$ and $\underline{\ell}''(\kappa)$, facilitates calculation of confidence intervals.

The constraints imposed on the κ parameter to positive numbers mean that it is more appropriate to work with the parameter in a transformed space,

$$\xi \equiv \log(\kappa) = g(\kappa).$$

Since the maximum likelihood estimator is asymptotically normally distributed and

$$g'(\kappa) = \kappa^{-1},$$

$$\hat{\xi} \sim N\left(\xi_{\text{true}}, \frac{1}{(\hat{\kappa})^2(-\ell''(\hat{\kappa}))}\right).$$

Thus for a 95% confidence interval we expect that

$$0.95 \approx P\left(\hat{\xi} - \frac{1.96}{\sqrt{(\hat{\kappa})^2(-\ell''(\hat{\kappa}))}} < \xi_{\text{true}} < \hat{\xi} + \frac{1.96}{\sqrt{(\hat{\kappa})^2(-\ell''(\hat{\kappa}))}}\right).$$

Setting

$$\xi_{\text{low}} = \log(\hat{\kappa}) - \frac{1.96}{\sqrt{(\hat{\kappa})^2(-\ell''(\hat{\kappa}))}} \quad \text{and} \quad \xi_{\text{upp}} = \log(\hat{\kappa}) + \frac{1.96}{\sqrt{(\hat{\kappa})^2(-\ell''(\hat{\kappa}))}},$$

then the lower and upper 95% confidence intervals for parameter κ are given by

$$\kappa_{\text{low}} = \exp(\xi_{\text{low}}) \quad \text{and} \quad \kappa_{\text{upp}} = \exp(\xi_{\text{upp}}).$$

5.3.4.1 Derivative approximation

We consider only the quasi-Newton solutions to the derivatives required for both point estimates and confidence intervals. To find the maximum of the likelihood surface and to obtain the second derivative, we use both the Nelder-Mead and BFGS algorithms via the R function `optim()` in the “stats” package (R Core Team, 2019⁵⁶). As this is discussed clearly in Section [2.4.2](#) we do not repeat it.

5.4 General negative binomial model

We now consider the extension to more general models where any number of fixed and random effects can be specified. This general model also allows for unbalanced datasets of m -groups with n_i observations per group. The form of this model is

$$y_i | \mathbf{u}_i \stackrel{\text{ind}}{\sim} \text{NB}\left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), \kappa; y_{ij}\right), \quad \mathbf{u}_i \stackrel{\text{ind}}{\sim} N(\mathbf{0}_{d_{\mathbf{R}}}, \boldsymbol{\Sigma}),$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i \quad \text{and} \quad \kappa > 0,$$

where the notation follows the general one level model presented in Section 1.8. The log-likelihood can be expressed as

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa) = \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa),$$

where

$$\begin{aligned} \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa) &= \log \int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{NB} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), \kappa; y_{ij} \right) \right\} \\ &\quad \times |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right) d\mathbf{u}_i, \end{aligned}$$

and the best predictor of \mathbf{u}_i is

$$\text{BP}(\mathbf{u}_i) \equiv \frac{\int_{\mathbb{R}^{d\mathbf{R}}} \mathbf{u}_i \left\{ \prod_{j=1}^{n_i} \text{NB} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), \kappa; y_{ij} \right) \right\} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right) d\mathbf{u}_i}{\int_{\mathbb{R}^{d\mathbf{R}}} \left\{ \prod_{j=1}^{n_i} \text{NB} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), \kappa; y_{ij} \right) \right\} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right) d\mathbf{u}_i}. \quad (5.25)$$

Implementation of EP follows as per the previous models, with minor alterations to the algorithm to account for the count response variable.

We first explain likelihood approximation using EP in Subsection 5.4.1, before explaining computation of point estimates and confidence intervals in Section 5.4.2. Results of a simulation study are presented in Subsection 5.4.3. The same calculations in Section 3.3 can be used for best prediction of \mathbf{u}_i . We do not repeat them again and instead refer readers to the previous work.

5.4.1 Expectation propagation likelihood approximation

EP can be used to approximate the likelihood by updating and summing the natural parameter updates. We wish to find the optimal natural parameters η_0 , $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ which minimise $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$, where f_{UN} is defined by equation (3.4) and

$$f_{\text{input}}(\mathbf{x}) = \text{NB} \left(\exp(c_0 + \mathbf{c}_1^\top \mathbf{x}), \kappa; c_2 \right) \exp \left(\left(\boldsymbol{\eta}_1^{\text{input}} \right)^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x} \right), \quad (5.26)$$

where $\boldsymbol{\eta}_1^{\text{input}}$ is a $d \times 1$ vector, $\mathbf{H}_2^{\text{input}}$ is a $d \times d$ matrix, $c_0 = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}}$, $c_2 = y_{ij}$, $\mathbf{c}_1 = \mathbf{x}_{ij}^{\mathbf{R}}$ and $\mathbf{x} = \mathbf{u}_i$. As such, we seek an $\boldsymbol{\eta}^*$ to solve

$$\int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \text{NB}\left(\exp(c_0 + \mathbf{c}_1^\top \mathbf{x}), \kappa; c_2\right) \exp\left\{\left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^{\text{input}}\right\} d\mathbf{x} = \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \exp\left\{\left[\begin{array}{c} 1 \\ \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^\top) \end{array}\right]^\top \boldsymbol{\eta}^*\right\} d\mathbf{x} \quad (5.27)$$

where $\mathbf{x}^{\otimes k}$ is as defined in equation (1.2). Thus to obtain the required projection, we first obtain the optimal natural parameters $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$ to project onto the multivariate normal family as is presented in Result 26. As the integrals required to do so are not available in closed form solutions, we use the same properties as the expit and Poisson models to express the multivariate integrals as univariate integrals. These integrals are expressed in Section 2.1 of Kim & Wand (2018)³² as denoted by $C_b(p, q, r)$ and solved using quadrature in a more efficient manner.

Result 26. Given f_{input} follows the form of equation (5.26), the projection onto the multivariate normal family is given by

$$\text{proj}_N[f_{\text{input}}] = \exp\left(\mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta}_{-1}^* - A(\boldsymbol{\eta}_{-1}^*)\right) h(\mathbf{x}),$$

where

$$\boldsymbol{\eta}_{-1}^* \equiv K_{NB}(\boldsymbol{\eta}_{-1}^{\text{input}}; c_0, \mathbf{c}_1, c_2; \kappa),$$

$$\boldsymbol{\eta}_{-1}^{\text{input}} \equiv \begin{bmatrix} \boldsymbol{\eta}_1^{\text{input}} \\ \boldsymbol{\eta}_2^{\text{input}} \end{bmatrix}, \quad \boldsymbol{\eta}_{-1}^* \equiv \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix},$$

$K_{NB}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2; \kappa\right)$ is defined in Definition 27 and $\mathbf{T}(\mathbf{x})$ and $h(\mathbf{x})$ follow from Section 1.5.2.2.

Definition 27. For primary arguments \mathbf{a}_1 ($d \times 1$) and \mathbf{a}_2 ($\frac{1}{2}d(d+1) \times 1$) such that $\text{vec}^{-1}\left(-(\mathbf{D}_d^+)^\top \mathbf{a}_2\right)$ is symmetric and positive definite, and auxiliary arguments $c_0, \in \mathbb{R}$, $c_2, \in \mathbb{Z}_{\geq 0}$, $\kappa \in \mathbb{R}_{\geq 0}$ and \mathbf{c}_1 ($d \times 1$), the function $K_{NB} : H \rightarrow H$ is given by

$$K_{NB}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2, \kappa\right) \equiv \begin{bmatrix} \mathbf{R}_5^\top (\mathbf{a}_1 + r_3 \mathbf{c}_1) \\ \mathbf{D}_d^\top \text{vec}(\mathbf{R}_5^\top \mathbf{A}_2) \end{bmatrix}, \quad (5.28)$$

where

$$\begin{aligned} r_1 &= \mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{c}_1, & r_2 &= (\mathbf{c}_1^\top \mathbf{A}_2^{-1} \mathbf{a}_1 - 2c_0)r_7, & r_3 &= r_2 + 2 \mathcal{C}_{b,1:0}(r_6, r_7, \kappa)r_7, \\ r_4 &= 2(\mathcal{C}_{b,1:0}(r_6, r_7, \kappa))^2 - \mathcal{C}_{b,2:0}(r_6, r_7, \kappa) - r_1)r_7^2, & \mathbf{R}_5 &= (\mathbf{A}_2 + r_4 \mathbf{c}_1 \mathbf{c}_1^\top)^{-1} \mathbf{A}_2, \\ r_6 &= 1 - r_2, & r_7 &= -r_1^{-1}, & \mathcal{C}_{b,1:0}(r_6, r_7, \kappa) &= \frac{\mathcal{C}_b(1, r_6, r_7, \kappa)}{\mathcal{C}_b(0, r_6, r_7, \kappa)}, \\ \mathcal{C}_{b,2:0}(r_6, r_7, \kappa) &= \frac{\mathcal{C}_b(2, r_6, r_7, \kappa)}{\mathcal{C}_b(0, r_6, r_7, \kappa)}, & b(x; \kappa) &= (c_2 + \kappa) \log(\exp(x) + \kappa) \end{aligned}$$

and $\mathcal{C}_b(k, r, q)$ follows from equation (5.42).

A proof of Definition 27 is given in Appendix 5.6.4. Using Result 26 we now obtain the normalising natural parameter η_0^* to find the projection onto unnormalised normal family.

5.4.1.1 Projection onto the unnormalised multivariate normal family

We now obtain the projection onto the unnormalised multivariate normal family. Recall the moment matching problem from equation (5.27) and Results 12 and 13. Then, we require

$$C_f = \int_{\mathbb{R}^d} f_{\text{input}}(x) dx = \mathcal{C}_b(0, r_6, r_7) Z_0 Z_1,$$

where

$$Z_1 = \exp\left(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) + \log \Gamma(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa)\right),$$

$Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi))$ and r_6 and r_7 are given in Definition 27. Analogous to previous arguments

$$\begin{aligned} \boldsymbol{\eta}_0^* &= \log \mathcal{C}_b(0, r_6, r_7) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) + \frac{1}{4} (\boldsymbol{\eta}_1^*)^\top (\mathbf{H}_2^*)^{-1} \boldsymbol{\eta}_1^* \\ &\quad - \frac{1}{4} (\boldsymbol{\eta}_1^{\text{input}})^\top (\mathbf{H}_2^{\text{input}})^{-1} \boldsymbol{\eta}_1^{\text{input}} + \frac{1}{2} \log(|\mathbf{H}_2^*|/|\mathbf{H}_2^{\text{input}}|) + \log \Gamma(c_2 + \kappa) \\ &\quad + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa). \end{aligned}$$

To find the normalising constant of the input function, we introduce a function given in Definition 28.

Definition 28. Consider first, primary arguments \mathbf{a}_1 and \mathbf{b}_1 and auxiliary argument \mathbf{c}_1 where all three are $d \times 1$. Next consider arguments \mathbf{a}_2 and \mathbf{b}_2 which are all $(\frac{1}{2}d(d+1) \times 1)$ such that both $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{a}_2)$ and $\text{vec}^{-1}(-(\mathbf{D}_d^+)^\top \mathbf{b}_2)$ are symmetric and positive definite. Finally note auxiliary scalar argument $c_0 \in \mathbb{R}$, $c_2 \in \mathbb{Z}_{\geq 0}$ and $\kappa \in \mathbb{R}_{\geq 0}$. Then the function $C_{NB} : H \times H \rightarrow \mathbb{R}$ is given by

$$\begin{aligned} C_{NB} \left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}; c_0, \mathbf{c}_1, c_2, \kappa \right) &\equiv \log \mathcal{C}_b(0, r_6, r_7, \kappa) + (r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi) \\ &\quad + \frac{1}{4} \mathbf{b}_1^\top \mathbf{B}_2^{-1} \mathbf{b}_1 - \frac{1}{4} \mathbf{a}_1^\top \mathbf{A}^{-1} \mathbf{a}_1 + \frac{1}{2} \log(|\mathbf{B}_2|/|\mathbf{A}_2|) \\ &\quad + \log \Gamma(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa), \end{aligned}$$

where $\mathbf{A}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{a}_2)$, $\mathbf{B}_2 \equiv \text{vec}^{-1}((\mathbf{D}_d^+)^\top \mathbf{b}_2)$, r_1 , r_2 , r_6 and r_7 follow from Definition 27.

In summary, the projection onto the unnormalised multivariate normal family is given by Result 27.

Result 27. For an unnormalised input function following the form of equation (5.26),

$$proj_{UN} [f_{input}] (x) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{x} \\ vec(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \begin{bmatrix} \eta_0^* \\ \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} \right\},$$

where

$$\begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix} = K_{NB} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}; c_0, \mathbf{c}_1, c_2, \kappa \right)$$

and

$$\eta_0^* = C_{NB} \left(\begin{bmatrix} \boldsymbol{\eta}_1^{input} \\ \boldsymbol{\eta}_2^{input} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\eta}_1^* \\ \boldsymbol{\eta}_2^* \end{bmatrix}; c_0, \mathbf{c}_1, c_2, \kappa \right).$$

5.4.1.2 Message passing formulation

Only minor changes to Algorithm 6 are required to account for the negative binomial case. Note the components of the sum of the log-likelihood function $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa)$ are

$$\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa) = \log \int_{\mathbb{R}^{d_R}} \left(\prod_{j=1}^{n_i} p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \right) p(\mathbf{u}_i; \boldsymbol{\Sigma}) d\mathbf{u}_i, \quad (5.29)$$

where

$$p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta}) \equiv \text{NB} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^F + \mathbf{u}_i^\top \mathbf{x}_{ij}^R), \kappa; y_{ij} \right)$$

and

$$p(\mathbf{u}_i; \boldsymbol{\Sigma}) \equiv |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i \right)$$

are the conditional density functions of each response given its random effect and the density function of that random effect respectively.

As the dependence structure of the product in equation (5.14) is the same as the probit model shown in Figure 3.1, we only need to update the messages depending on the conditional density functions of each response given its random effect. This involves

updating messages $m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}(\mathbf{u}_i)$ as in equation (3.12)

$$m_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}(\mathbf{u}_i) \leftarrow \frac{\text{proj}_{\text{UN}} \left[\text{NB} \left(\exp(c_0 + \mathbf{c}_{1_{ij}}^\top \mathbf{u}_i), \kappa; c_{2_{ij}} \right) \exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\} \right]}{\exp \left\{ \mathbf{u}_i^\top \boldsymbol{\eta}_1^\otimes + (\text{vech}(\mathbf{u}_i \mathbf{u}_i^\top))^\top \boldsymbol{\eta}_2^\otimes \right\}},$$

where we set the constant terms:

$$c_{0_{ij}} \leftarrow \boldsymbol{\beta}^\top \mathbf{x}_{ij}^\mathbf{F}; \quad \mathbf{c}_{1_{ij}} \leftarrow \mathbf{x}_{ij}^\mathbf{R}; \quad c_{2_{ij}} \leftarrow y_{ij}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i.$$

Following Section 3.1.2, we see that this task reduces down to adjusting the calculation of the optimal natural parameters $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ in equation (3.18) for the Poisson model, where the linear and quadratic coefficient updates in equation (3.19) given by K_{probit} are changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i})_{1:2} \leftarrow K_{\text{NB}}(\boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ij}}, \mathbf{c}_{1_{ij}}, c_{2_{ij}}, \kappa) - \boldsymbol{\eta}_{1:2}^\otimes \quad (5.30)$$

as per Definition 27, and where the constant coefficient update in equation (3.20) given by C_{probit} is changed to

$$(\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i})_0 \leftarrow C_{\text{NB}}(\boldsymbol{\eta}_{1:2}^\otimes, (\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i})_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ij}}, \mathbf{c}_{1_{ij}}, c_{2_{ij}}, \kappa)$$

as per Definition 28.

Barring the aforementioned changes and the initialisation of $\boldsymbol{\eta}_{p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta})\rightarrow\mathbf{u}_i}$ with the ones discussed in Section 5.4.1.3, the full algorithm for the approximation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ follows Algorithm 6 exactly.

5.4.1.3 Starting values for the negative binomial case

We now derive starting values for the negative binomial case using the same principles shown in the other cases. Let

$$\log p(y_{ij}|\mathbf{u}_i;\boldsymbol{\beta}) = f(a_{ij}) + \log \Gamma(y_{ij} + \kappa) - \log \Gamma(y_{ij} + 1) - \log \Gamma(\kappa) + \kappa \log \kappa,$$

where $a_{ij} = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^\mathbf{F} + \mathbf{u}_i^\top \mathbf{x}_{ij}^\mathbf{R}$ and $f(x)$ is defined as in equation (5.23). Then

$$(a_{ij} - \hat{a}_{ij}) = (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^\mathbf{R}$$

where $\hat{a}_{ij} = \boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \hat{\mathbf{u}}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}$ and $\hat{\mathbf{u}}$ is an approximation of \mathbf{u}_i . A Taylor series expansion of $f(a_{ij})$ evaluated at \hat{a}_{ij} leads to

$$\begin{aligned} f(a_{ij}) &= f(\hat{a}_{ij}) + f'(\hat{a}_{ij})(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}} + \frac{1}{2} f''(\hat{a}_{ij})((\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{x}_{ij}^{\mathbf{R}})^2 + \dots \\ &= \begin{bmatrix} 1 \\ \mathbf{u}_i - \hat{\mathbf{u}}_i \\ \text{vech}((\mathbf{u}_i - \hat{\mathbf{u}}_i)(\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top) \end{bmatrix}^\top \check{\boldsymbol{\eta}}_{ij} + \dots, \end{aligned}$$

where

$$\check{\boldsymbol{\eta}}_{ij} = \begin{bmatrix} f(\hat{a}_{ij}) \\ f'(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} \\ \frac{1}{2} f''(\hat{a}_{ij}) \mathbf{D}_{d^{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix}.$$

Quadratic approximation to $\log p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})$ based on Taylor expansion about $\hat{\mathbf{u}}_i$ is $\log \check{p}(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})$ as per equation (3.22). By following the same logic of previous cases it is easy to show

$$\boldsymbol{\eta}_{p(y_{ij} | \mathbf{u}_i; \boldsymbol{\beta})}^{\text{start}} = \begin{bmatrix} \eta_0^{\text{start}} \\ f'(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} - f''(\hat{a}_{ij}) \mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_i \\ \frac{1}{2} f''(\hat{a}_{ij}) \mathbf{D}_{d^{\mathbf{R}}}^\top \text{vec}(\mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top) \end{bmatrix},$$

where

$$\eta_0^{\text{start}} = f(\hat{a}_{ij}) - f'(\hat{a}_{ij}) (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_{ij} + \frac{1}{2} f''(\hat{a}_{ij}) \hat{\mathbf{u}}_{ij}^\top \mathbf{x}_{ij}^{\mathbf{R}} (\mathbf{x}_{ij}^{\mathbf{R}})^\top \hat{\mathbf{u}}_{ij}.$$

5.4.2 Computation of point estimates and confidence intervals

We denote the maximum likelihood approximation using EP as $\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Sigma}}, \hat{\kappa} = \underset{\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa}{\text{argmax}} \ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa)$.

As before, to find their stationary points, we use the Nelder-Mead and BFGS algorithms for optimisation via the R function `optim()` in the package “stats” (R Core Team, 2019^[56]). As in the κ -only model, additional care must be taken to ensure optimisation of the parameter κ is on an unconstrained space. We do this using a log transform as in the univariate case i.e. $\xi = \log(\kappa)$. For clarity we now show the changes required from the steps given in Section 3.2:

1. Steps (a) - (b) for converting $\boldsymbol{\Sigma}$ to the unconstrained space $\boldsymbol{\theta}$ follows from before.

We must also now convert κ to the unconstrained space ξ .

- (c) Convert κ to the unconstrained space ξ

$$\xi = \log(\kappa).$$

2. A quasi-Newton optimisation method can now be used to obtain the maximum likelihood estimate of $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\xi})$,

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\xi}) = \underset{\boldsymbol{\theta} \in \mathbb{R}^{d^{\boldsymbol{\Sigma}}(d^{\boldsymbol{\Sigma}}+1)/2}}{\operatorname{argmax}} \ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \xi).$$

We suggest conducting an initial search via the Nelder-Mead method, with refinements by BFGS algorithm. Both can be implemented via the `optim()` R function in the “stats” package (R Core Team, 2019^[56]).

3. The step converting $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}})$ to $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}})$ follows from before.
4. Now obtain the Hessian matrix $H\ell(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\xi})$ at the maximum $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\xi})$ using the quasi-Newton method BFGS, which as before can be implemented via `optim()`. Although we need values of $(\boldsymbol{\beta}, \boldsymbol{\omega}, \xi)$ to be returned in the Hessian, the constraints on these parameters mean the Hessian should still be calculated on the $(\boldsymbol{\beta}, \boldsymbol{\theta}, \xi)$ space, which is computed as before.
5. Form $100(1 - \alpha)\%$ confidence intervals for the entries of $(\boldsymbol{\beta}, \boldsymbol{\omega})$ using

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\omega}} \\ \hat{\xi} \end{bmatrix} \pm \Phi^{-1} \left(1 - \frac{1}{2}\alpha \right) \sqrt{-\operatorname{diag} \left\{ \left(\mathbf{H}\ell(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\xi}) \right)^{-1} \right\}}.$$

6. Back transform the confidence interval limits for $\boldsymbol{\omega}$ corresponding to the standard deviation and correlation parameters as before and back transform ξ to κ by exponentiating it.

5.4.3 Simulation study

We now provide a simulation study repeated 1000 times comparing EP to a quadrature approach and Laplace approximation. The simulation study in this chapter does not assess the speed component.

Datasets were generated with true parameter values:

$$\beta_{\text{true}} = [0.38, 0.93]^\top, \quad \sigma_{\text{true}}^2 = 0.25 \quad \text{and} \quad \kappa_{\text{true}} = 15. \quad (5.31)$$

There were 100 groups generated in the data with each group containing 5 measurements (i.e. $m = 100$, $n = 5$). The $\mathbf{x}_{ij}^{\mathbf{F}}$ and $\mathbf{x}_{ij}^{\mathbf{R}}$ vectors were of the form

$$\mathbf{x}_{ij}^{\mathbf{F}} = [1, x_{1ij}]^\top \quad \text{and} \quad \mathbf{x}_{ij}^{\mathbf{R}} = 1,$$

where $x_{k,ij}$ was generated independently from a uniform distribution on the unit interval. The tolerance of error for the EP scheme was set to 10^{-5} and 1000 iterations were allowed for optimisation. We compare the EP approach to the Laplace approximation and 100 point adaptive Gauss-Hermite quadrature. Both alternative approaches were implemented via the R function `g1mer()` from the package “lme4” (Bates, et al., 2018⁵).

Point estimates and 95% confidence intervals for each interpretable model parameter of the study are presented in Figure 5.4, where the number in the upper-right hand corner of each panel are the empirical coverage values based on all 1000 replicates. Only 20 randomly chosen replicates from each method are shown in the panels for ease of viewing. Laplace approximation and EP are shown in black, super imposed on adaptive Gauss-Hermite quadrature approach shown in grey.

None of the three methods provided over 95% empirical coverage for the fixed intercept. Notably they were within 0.3% of the highest coverage, which was EP at 94.4%. Although for the fixed slope all methods had coverage above 95%, the quadrature approach delivered the highest coverage of 95.9%, 0.4% above the coverage of EP and 0.7% above Laplace approximations. EP provided the best estimates of 94.3% for the random intercept parameter by some margin. Surprisingly, quadrature provided coverage than Laplace approximation for this parameter, with 81.3% opposed to 92.8%.

Neither Laplace approximation nor Gauss-Hermite quadrature support confidence interval calculation for the shape parameter via the R function `g1mer()`. Differences between the point estimates of Laplace approximation and adaptive Gauss-Hermite quadrature exist for small κ values, although for larger values they are minimal. Point estimates of EP differ from adaptive Gauss-Hermite quadrature across the range of values, however match with Laplace approximation values for small values of κ . Although confidence intervals via EP provides reasonable coverage of 94.4%, the large width of intervals may not lend itself well to meaningful inference. It is left as an open topic for future research.

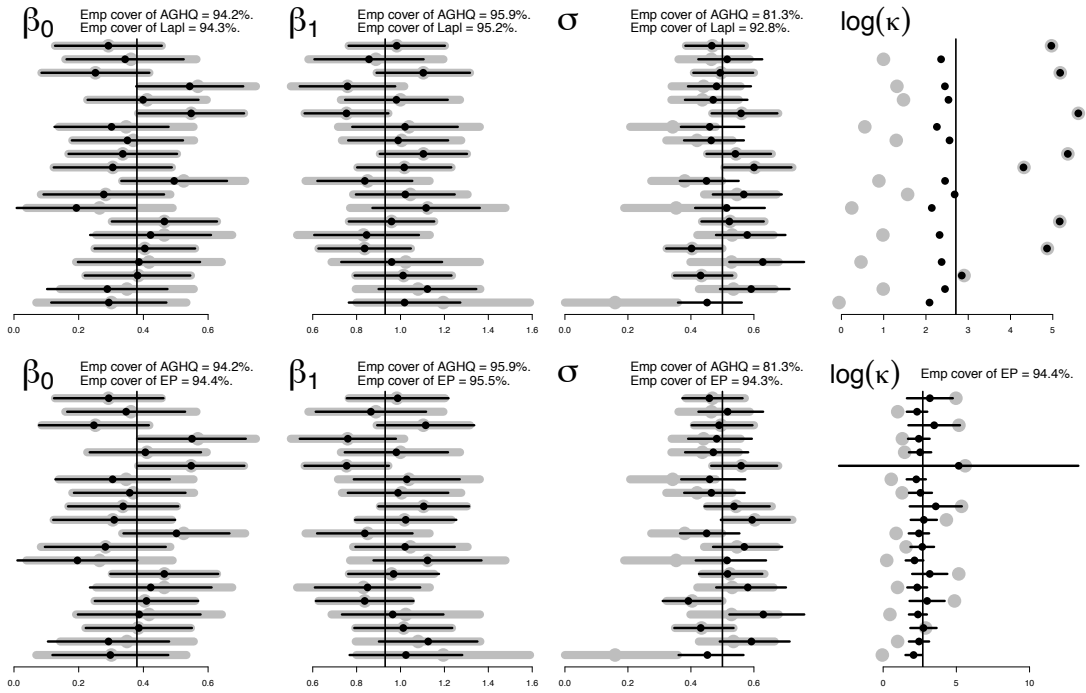


Figure 5.4: Plots of confidence interval coverage for the univariate model with true parameter values from equation (5.31) for models fit with adaptive Gauss-Hermite quadrature, Laplace approximation and EP. The horizontal lines are the confidence intervals for 20 randomly chosen replications of the simulation study, the solid circular points indicate the point estimates and the vertical lines indicate true parameter values. The Laplace and EP are shown in black, superimposed over the grey adaptive Gauss-Hermite quadrature estimates. The percentage given in the top right-hand corner of each panel is the empirical coverage over all 1000 replications.

5.5 Varying dispersion negative binomial model

We give an indication of how the proposed general negative binomial model can be extended to the case of varying dispersion, where each group has its own shape parameter. Rather than estimating a unique κ for each group, we instead estimate an error term ν_i for each group and multiply κ by the exponent of this term. The form of this model is

$$y_{ij} | \mathbf{u}_i \stackrel{\text{ind}}{\sim} \text{NB} \left(\exp \left(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_{ij}^\top \mathbf{x}_{ij}^{\mathbf{R}} \right), \kappa \exp(\nu_i) \right), \quad \mathbf{u}_i \stackrel{\text{ind}}{\sim} N(\mathbf{0}_{d\mathbf{R}}, \boldsymbol{\Sigma}), \quad \nu_i \stackrel{\text{ind}}{\sim} N(0, \sigma_\nu^2),$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i \quad \kappa > 0, \quad \text{and} \quad \nu_i \in \mathbb{R},$$

where the notation follows the general one level model presented in Section 1.8. The log-likelihood can be expressed as

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa, \sigma_v^2) = \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa, \sigma_v^2),$$

where

$$\begin{aligned} \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \kappa, \sigma_v^2) = & \log \int_{\mathbb{R}^{d_{\mathbf{R}}}} \int_{-\infty}^{\infty} \left\{ \prod_{j=1}^{n_i} \text{NB} \left(\exp(\boldsymbol{\beta}^\top \mathbf{x}_{ij}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ij}^{\mathbf{R}}), \kappa; y_{ij} \right) \right\} \\ & \times |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i\right) (2\pi\sigma_v^2)^{-1/2} \exp\left(-\frac{\nu_i^2}{2\sigma_v^2}\right) d\nu_i d\mathbf{u}_i. \end{aligned} \quad (5.32)$$

Although it seems that the extension of our methodology to this model is complicated, the likelihood as in equation (5.32) can be obtained by message passing on the updated factor graph in Figure 5.5 in a similar way to the general negative binomial model. We leave further development of this model for future research.

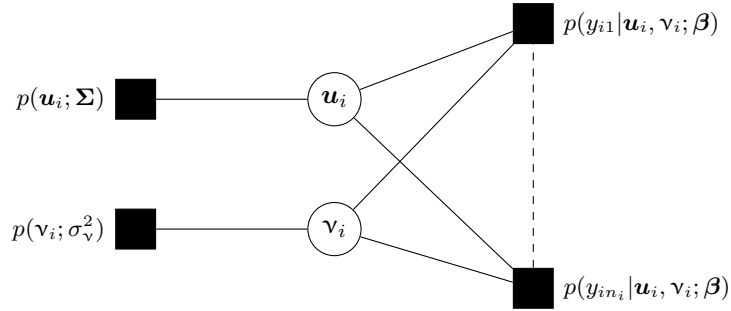


Figure 5.5: Factor graph representation of the product structure of the integrand in equation (5.32).

5.6 Appendix

5.6.1 Proof of Definition 21

From this point on we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. The k th moment of the input function in equation (5.2) is given by

$$\begin{aligned} \mathcal{M}_k &= \int_{-\infty}^{\infty} x^k \Psi(c_0 + c_1 x; c_2) (2\pi)^{-1/2} \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} dx \\ &\quad \times (2\pi)^{1/2} \exp(A(\boldsymbol{\eta})) \Gamma(c_2 + 1)^{-1}, \end{aligned}$$

where $\Psi(x; c) = \exp(cx - b(x))$ and $b(x) = \exp(x)$. Using the inverse map of the natural parameters in equation (1.14) and implementing the change of variable $x = \mu + \sigma u$,

$$\mathcal{M}_k = \int_{-\infty}^{\infty} (\mu + \sigma u)^k \Psi(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{1}{2} \log(2\pi) - \log \Gamma(c_2 + 1))$. Each of the required moments can be obtained via univariate quadrature. For numerical stability we rearrange the integral that arises in the form presented by Kim & Wand (2017),^[32]

$$\mathcal{C}_b(k, r, q) = \int_{-\infty}^{\infty} x^k \exp(rx - qx^2 - b(x)) dx, \quad (5.33)$$

where $b(x) = \exp(x)$. Using simple algebraic manipulations, it is easy to show the integrals required for the calculation of each moment can be expressed as

$$\begin{aligned} &\int_{-\infty}^{\infty} u^k \Psi(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du \\ &= \int_{-\infty}^{\infty} ((r_2 + 2r_7 x)(2r_7)^{-1/2})^k \exp(r_6 x - r_7 x^2 - b(x)) dx Z_0, \end{aligned}$$

where

$$\begin{aligned} r_1 &= -2\sigma^2 c_1^2, \quad r_2 = 2(c_0 + c_1 \mu) r_1^{-1}, \quad r_6 = c_2 - r_2, \quad r_7 = -r_1^{-1}, \\ b(x) &= \exp(x) \quad \text{and} \quad Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi)). \end{aligned}$$

The explicit forms of the integrals required follow easily,

$$\begin{aligned} Z_0^{-1} \int_{-\infty}^{\infty} \Psi(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} \int_{-\infty}^{\infty} u \Psi(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du \\ &= (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) (2r_7)^{-1/2}, \\ Z_0^{-1} \int_{-\infty}^{\infty} u^2 \Psi(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du \\ &= 2(r_2 \mathcal{C}_b(1, r_6, r_7) + r_7 \mathcal{C}_b(2, r_6, r_7)) - \frac{r_2^2 r_1}{2} \mathcal{C}_b(0, r_6, r_7). \end{aligned}$$

Implementing these forms, each of the moments are

$$\begin{aligned} Z_0^{-1} Z_1^{-1} \mathcal{M}_0 &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_1 &= \mu \mathcal{C}_b(0, r_6, r_7) + \sigma (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) (2r_7)^{-1/2}, \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_2 &= (\mu^2 + \sigma^2) \mathcal{C}_b(0, r_6, r_7) + 2\mu c_1 \sigma^2 (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) \\ &\quad + \sigma^4 c_1^2 (4r_7^2 \mathcal{C}_b(2, r_6, r_7) + 4r_7 r_2 \mathcal{C}_b(1, r_6, r_7) + (r_2^2 - 2r_7) \mathcal{C}_b(0, r_6, r_7)). \end{aligned}$$

Letting $\mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$ and $\mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$, the optimal mean parameter for the projection is

$$\mu^* = E(x) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \mu + \sigma^2 c_1 (r_2 + 2r_7 \mathcal{C}_{b,1:0}(r_6, r_7)),$$

and the optimal variance parameter is

$$\begin{aligned} (\sigma^2)^* &= E(x^2) - E(x)^2 = \frac{\mathcal{M}_2}{\mathcal{M}_0} - \left(\frac{\mathcal{M}_1}{\mathcal{M}_0} \right)^2 \\ &= \sigma^2 + (2r_7)^2 \sigma^4 c_1^2 (\mathcal{C}_{b,2:0}(r_6, r_7) - \mathcal{C}_{b,1:0}(r_6, r_7)^2 + r_1/2). \end{aligned}$$

By converting back to natural optimal and input parameters we arrive at Definition [17](#).

5.6.2 Proof of Definition [23](#)

Using simple algebraic manipulations based on Lemma [2](#), we arrive at Lemma [5](#):

Lemma 5. *For integrals of the forms listed below, the corresponding solutions exist:*

$$\int_{\mathbb{R}^d} \Psi(a + \mathbf{b}^\top \mathbf{x}; c) \phi(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} \Psi(a + \|\mathbf{b}\|x; c) \phi(x) dx, \quad (5.34)$$

$$\int_{\mathbb{R}^d} \mathbf{x} \Psi(a + \mathbf{b}^\top \mathbf{x}; c) \phi(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{b}}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} x \Psi(a + \|\mathbf{b}\|x; c) \phi(x) dx, \quad (5.35)$$

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^\top \Psi(a + \mathbf{b}^\top \mathbf{x}; c) \phi(\mathbf{x}) d\mathbf{x} &= I_d \int_{-\infty}^{\infty} \Psi(a + \|\mathbf{b}\|x; c) \phi(x) dx \\ &\quad + \frac{\mathbf{b} \mathbf{b}^\top}{\mathbf{b}^\top \mathbf{b}} \left(\int_{-\infty}^{\infty} x^2 \Psi(a + \|\mathbf{b}\|x; c) \phi(x) dx - \int_{-\infty}^{\infty} \Psi(a + \|\mathbf{b}\|x; c) \phi(x) dx \right), \end{aligned} \quad (5.36)$$

where $a \in \mathbb{R}$, \mathbf{b} is a $d \times 1$ vector, $c \in \mathbb{Z}^+$, $\Psi(x; c) = \exp(cx - b(x))$ and $b(x) = \exp(x)$.

We wish to obtain the projection of an input function following the form of equation [\(5.11\)](#) onto the multivariate normal family. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. Using the $\mathbf{x}^{\otimes k}$ notation as described in equation [\(1.2\)](#), the k th moment of the input function is given by

$$\begin{aligned} \mathcal{M}_k &= \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Psi(c_0 + \mathbf{c}_1^\top \mathbf{x}; c_2) (2\pi)^{-d/2} \exp \left\{ \left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x} \mathbf{x}^\top) \end{array} \right]^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} d\mathbf{x} \\ &\quad \times (2\pi)^{d/2} \exp(A(\boldsymbol{\eta})) \Gamma(c_2 + 1)^{-1}. \end{aligned}$$

Using the inverse map of the natural parameters in equation (1.14) and implementing a change of variable, where $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\mathbf{u}$,

$$\mathcal{M}_k = \int_{\mathbb{R}^d} (\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\mathbf{u})^{\otimes k} \Psi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + (\boldsymbol{\Sigma}^{1/2}\mathbf{c}_1)^\top \mathbf{u}; c_2) \phi_{\mathbf{I}}(\mathbf{u}) d\mathbf{u} Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) - \log \Gamma(c_2 + 1))$. We can then use Lemma 5 to obtain each of the required moments via univariate quadrature. For numerical stability we rearrange the integral that arises in the form presented by Kim & Wand (2017)³²

$$\mathcal{C}_b(k, r, q) = \int_{-\infty}^{\infty} x^k \exp(rx - qx^2 - b(x)) dx, \quad (5.37)$$

where $b(x) = \exp(x)$. Using simple algebraic manipulations, it is easy to show

$$\begin{aligned} & \int_{-\infty}^{\infty} u^k \Psi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2}\mathbf{c}_1\|u; c_2) \phi(u) du \\ &= \int_{-\infty}^{\infty} ((r_2 + 2r_7x)(2r_7)^{-1/2})^k \exp(r_6x - r_7x - b(x)) dx Z_0, \end{aligned}$$

where

$$\begin{aligned} r_1 &= -2\mathbf{c}_1^\top \boldsymbol{\Sigma} \mathbf{c}_1, \quad r_2 = 2(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu})r_1^{-1}, \quad r_6 = c_2 - r_2, \quad r_7 = -r_1^{-1}, \\ b(x) &= \exp(x) \quad \text{and} \quad Z_0 = \exp((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi)). \end{aligned}$$

It follows,

$$\begin{aligned} Z_0^{-1} \int_{-\infty}^{\infty} \Psi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2}\mathbf{c}_1\|u; c_2) \phi(u) du &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} \int_{-\infty}^{\infty} u \Psi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2}\mathbf{c}_1\|u; c_2) \phi(u) du \\ &= (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7))(2r_7)^{-1/2}, \\ Z_0^{-1} \int_{-\infty}^{\infty} u^2 \Psi(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2}\mathbf{c}_1\|u; c_2) \phi(u) du \\ &= 2(r_2 \mathcal{C}_b(1, r_6, r_7) + r_7 \mathcal{C}_b(2, r_6, r_7)) - \frac{r_2^2 r_1}{2} \mathcal{C}_b(0, r_6, r_7). \end{aligned}$$

Implementing these forms, each of the moments can be derived in a similar way to the probit cases. For the zeroth moment where $k = 0$,

$$\begin{aligned} Z_0^{-1} Z_1^{-1} \mathcal{M}_0 &= \mathcal{C}_b(0, r_6, r_7), \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_1 &= \boldsymbol{\mu} \mathcal{C}_b(0, r_6, r_7) + \boldsymbol{\Sigma}^{1/2} (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) (2r_7)^{-1/2}, \\ Z_0^{-1} Z_1^{-1} \mathcal{M}_2 &= (\boldsymbol{\mu} \boldsymbol{\mu}^\top + \boldsymbol{\Sigma}) \mathcal{C}_b(0, r_6, r_7) \\ &\quad + (\boldsymbol{\mu} \mathbf{c}_1^\top \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{c}_1 \boldsymbol{\mu}^\top) (r_2 \mathcal{C}_b(0, r_6, r_7) + 2r_7 \mathcal{C}_b(1, r_6, r_7)) \\ &\quad + \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma} (4r_7^2 \mathcal{C}_b(2, r_6, r_7) + 4r_7 r_2 \mathcal{C}_b(1, r_6, r_7) + (r_2^2 - 2r_7) \mathcal{C}_b(0, r_6, r_7)). \end{aligned}$$

First, let $\mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$ and $\mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$. Then the optimal mean parameter for the projection is

$$\boldsymbol{\mu}^* = E(\mathbf{x}) = \frac{\mathcal{M}_1}{\mathcal{M}_0} = \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{c}_1 (r_2 + 2r_7 \mathcal{C}_{b,1:0}(r_6, r_7)),$$

and the optimal variance parameter is

$$\begin{aligned} \boldsymbol{\Sigma}^* &= E(\mathbf{x}\mathbf{x}^\top) - E(\mathbf{x})E(\mathbf{x})^\top = \frac{\mathcal{M}_2}{\mathcal{M}_0} - \frac{\mathcal{M}_1}{\mathcal{M}_0} \left(\frac{\mathcal{M}_1}{\mathcal{M}_0} \right)^\top \\ &= \boldsymbol{\Sigma} + (2r_7)^2 \boldsymbol{\Sigma} \mathbf{c}_1 \mathbf{c}_1^\top \boldsymbol{\Sigma} (\mathcal{C}_{b,2:0}(r_6, r_7) - \mathcal{C}_{b,1:0}(r_6, r_7)^2 + r_1/2), \end{aligned}$$

where $\mathcal{C}_{b,1:0}(r_6, r_7) = \frac{\mathcal{C}_b(1, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$ and $\mathcal{C}_{b,2:0}(r_6, r_7) = \frac{\mathcal{C}_b(2, r_6, r_7)}{\mathcal{C}_b(0, r_6, r_7)}$. By converting back to natural optimal and input parameters we arrive at Definition [17](#)

5.6.3 Proof of Definition [25](#)

We wish to obtain the projection of an input function following the form of equation [5.18](#) onto the univariate normal family. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. Using the x^k notation as described in equation [1.2](#), the k th moment of the input function is given by

$$\begin{aligned} \mathcal{M}_k &= \int_{\mathbb{R}} x^k \Upsilon(c_2; c_0 + c_1^\top x, \kappa) (2\pi)^{-1/2} \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} dx \\ &\quad \times \frac{\Gamma(c_2 + \kappa) \kappa^\kappa}{\Gamma(c_2 + 1) \Gamma(\kappa)} (2\pi)^{1/2} \exp(A(\boldsymbol{\eta})), \end{aligned}$$

where $\Upsilon(c; x; \kappa) = \exp(cx - b(x; c, \kappa))$ and $b(c; x; \kappa) = (c + \kappa) \log(\exp(x) + \kappa)$. Using the inverse map of the natural parameters in equation [1.14](#) and implementing a change of variable $x = \mu + \sigma u$

$$\mathcal{M}_k = \int_{\mathbb{R}} (\mu + \sigma u)^k \Upsilon(c_2; c_0 + c_1 \mu + (\sigma c_1)^\top u; \kappa) \phi(u) du Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{1}{2} \log(2\pi) + \log(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa))$. We can then use Lemma [6](#) to obtain each of the required moments via univariate quadrature. For numerical stability we rearrange the integral that arise in the form presented by Kim & Wand (2017)[32](#)

$$\mathcal{C}_b(k, r, q) = \int_{-\infty}^{\infty} x^k \exp(rx - qx^2 - b(x; c_2, \kappa)) dx. \quad (5.38)$$

Using simple algebraic manipulations, it is easy to show

$$\begin{aligned} &\int_{-\infty}^{\infty} u^k \Upsilon(c_0 + c_1 \mu + \sigma c_1 u; c_2) \phi(u) du \\ &= \int_{-\infty}^{\infty} ((r_2 + 2r_7 x) (2r_7)^{-1/2})^k \exp(r_6 x - r_7 x - b(x)) dx Z_0, \end{aligned}$$

where,

$$r_1 = -2c_1^\top \sigma c_1, \quad r_2 = 2(c_0 + c_1^\top \mu)r_1^{-1}, \quad r_6 = c_2 - r_2, \quad r_7 = -r_1^{-1},$$

$$b(x) = \exp(x) \quad \text{and} \quad Z_0 = \exp\left((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi)\right).$$

Using simple algebraic manipulations analogous to the Poisson case in Section 5.6.1 it is easy to arrive at Definition 25.

5.6.4 Proof of Definition 27

Using simple algebraic manipulations based on Lemma 2, we arrive at Lemma 6:

Lemma 6. *For integrals of the forms listed below, the corresponding solutions exist:*

$$\int_{\mathbb{R}^d} \Upsilon(c, a + \mathbf{b}^\top \mathbf{x}, \kappa) \phi(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} \Upsilon(c, a + \|\mathbf{b}\|x, \kappa) \phi(x) dx, \quad (5.39)$$

$$\int_{\mathbb{R}^d} \mathbf{x} \Upsilon(c, a + \mathbf{b}^\top \mathbf{x}, \kappa) \phi(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{b}}{\|\mathbf{b}\|} \int_{-\infty}^{\infty} x \Upsilon(c, a + \|\mathbf{b}\|x, \kappa) \phi(x) dx, \quad (5.40)$$

$$\int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^\top \Upsilon(c, a + \mathbf{b}^\top \mathbf{x}, \kappa) \phi(\mathbf{x}) d\mathbf{x} = I_d \int_{-\infty}^{\infty} \Upsilon(c, a + \|\mathbf{b}\|x, \kappa) \phi(x) dx$$

$$+ \frac{\mathbf{b} \mathbf{b}^\top}{\mathbf{b}^\top \mathbf{b}} \left(\int_{-\infty}^{\infty} x^2 \Upsilon(c, a + \|\mathbf{b}\|x, \kappa) \phi(x) dx - \int_{-\infty}^{\infty} \Upsilon(c, a + \|\mathbf{b}\|x, \kappa) \phi(x) dx \right), \quad (5.41)$$

where $a \in \mathbb{R}$, \mathbf{b} is a $d \times 1$ vector, $c \in \mathbb{Z}^+$, $\Upsilon(c, x, \kappa) = \exp(cx - b(c, x, \kappa))$ and $b(c, x, \kappa) = (c + \kappa) \log(\exp(x) + \kappa)$.

We wish to obtain the projection of an input function following the form of equation (5.11) onto the multivariate normal family. Note in the interest of brevity we represent the input parameter $\boldsymbol{\eta}^{\text{input}}$ as $\boldsymbol{\eta}$. Using the $\mathbf{x}^{\otimes k}$ notation as described in equation (1.2), the k th moment of the input function is given by

$$\mathcal{M}_k = \int_{\mathbb{R}^d} \mathbf{x}^{\otimes k} \Upsilon(c_2, c_0 + \mathbf{c}_1^\top \mathbf{x}, \kappa) (2\pi)^{-d/2} \exp \left\{ \left[\begin{array}{c} \mathbf{x} \\ \text{vech}(\mathbf{x} \mathbf{x}^\top) \end{array} \right]^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) \right\} dx$$

$$\times \frac{\Gamma(c_2 + \kappa) \kappa^\kappa}{\Gamma(c_2 + 1) \Gamma(\kappa)} (2\pi)^{d/2} \exp(A(\boldsymbol{\eta})).$$

Using the inverse map of the natural parameters in equation (1.14) and implementing a change of variable where $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u}$,

$$\mathcal{M}_k = \int_{\mathbb{R}^d} (\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{u})^{\otimes k} \Upsilon(c_2, c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + (\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1)^\top \mathbf{u}, \kappa) \phi_{\mathbf{I}}(\mathbf{u}) d\mathbf{u} Z_1,$$

where $Z_1 = \exp(A(\boldsymbol{\eta}) + \frac{d}{2} \log(2\pi) + \log(c_2 + \kappa) + \kappa \log \kappa - \log \Gamma(c_2 + 1) - \log \Gamma(\kappa))$. We can then Lemma 6 to obtain each of the required moment via univariate quadrature.

For numerical stability we rearrange the integral that arises in the form presented by Kim & Wand (2017),⁵²

$$\mathcal{C}_b(k, r, q) = \int_{-\infty}^{\infty} x^k \exp(rx - qx^2 - b(x)) dx, \quad (5.42)$$

where $b(x, \kappa) = (c_2 + \kappa) \log(\exp(x) + \kappa)$. Using simple algebraic manipulations, it is easy to show

$$\begin{aligned} & \int_{-\infty}^{\infty} u^k \Upsilon(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu} + \|\boldsymbol{\Sigma}^{1/2} \mathbf{c}_1\| u, c_2) \phi(u) du \\ &= \int_{-\infty}^{\infty} ((r_2 + 2r_7 x)(2r_7)^{-1/2})^k \exp(r_6 x - r_7 x - b(x)) dx Z_0, \end{aligned}$$

where

$$r_1 = -2\mathbf{c}_1^\top \boldsymbol{\Sigma} \mathbf{c}_1, \quad r_2 = 2(c_0 + \mathbf{c}_1^\top \boldsymbol{\mu}) r_1^{-1}, \quad r_6 = c_2 - r_2, \quad r_7 = -r_1^{-1},$$

$$b(x) = \exp(x) \quad \text{and} \quad Z_0 = \exp\left((r_2/2)^2 r_1 + \frac{1}{2} \log(r_7/\pi)\right).$$

By following the Poisson case in Section 5.6.2 it is easy to arrive at Definition 27.

Chapter 6

Expectation propagation for two level and crossed random effects probit models

Having explored models for both count and binary response data with one level of nesting, we now attempt to handle models for crossed random effects and two level structures. We work with probit models since they have closed form solutions as in Chapter 3 and show how the same key results from Chapter 3 are applicable to the higher level models in this chapter. We start by applying our work to the crossed random effects model in Section 6.1, before explaining how it can also be implemented for two level models in Section 6.2

6.1 The general probit crossed mixed model

We first extend our methodology from the one level probit case to crossed random effects GLMMs. We aim to find approximations of the maximum likelihood estimates for parameters β , Σ and Σ' with 95% confidence intervals.

The crossed random effects model specification is given in equation (6.1),

$$\begin{aligned} y_{ii'j} | \mathbf{u}_i, \mathbf{u}_{i'} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left\{ \Phi \left(\beta^\top \mathbf{x}_{ii'j}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ii'j}^{\mathbf{R}} + (\mathbf{u}_{i'})^\top \mathbf{x}_{ii'j}^{\mathbf{R}'} \right) \right\}, \\ \mathbf{u}_i &\stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d_{\mathbf{R}}}, \Sigma) \text{ independently of } \mathbf{u}_{i'} \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d_{\mathbf{R}'}}', \Sigma'), \\ 1 \leq i \leq m, \quad 1 \leq i' \leq m', \quad 1 \leq j \leq n_{ii'}, \end{aligned} \tag{6.1}$$

where the notation follows that of the crossed random effects model in Section [1.8](#).

The form of the log-likelihood is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') = \log \int_{\mathbb{R}^{m d^{\mathbf{R}} + m' d^{\mathbf{R}'}}} \left(\prod_{(i,i'): n_{ii'} > 0} \prod_{j=1}^{n_{ii'}} p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \right) p(\mathbf{u}, \mathbf{u}'; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') d \begin{bmatrix} \mathbf{u} \\ \mathbf{u}' \end{bmatrix}, \quad (6.2)$$

where

$$p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) = \Phi \left\{ (2y_{ii'j} - 1) \left(\boldsymbol{\beta}^\top \mathbf{x}_{ii'j}^{\mathbf{F}} + \mathbf{u}_i^\top \mathbf{x}_{ii'j}^{\mathbf{R}} + (\mathbf{u}'_{i'})^\top \mathbf{x}_{ii'j}^{\mathbf{R}'} \right) \right\},$$

$$p(\mathbf{u}, \mathbf{u}'; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \text{ is } \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_{m d^{\mathbf{R}}} \\ \mathbf{0}_{m' d^{\mathbf{R}'}} \end{bmatrix}, \begin{bmatrix} \mathbf{I}_m \otimes \boldsymbol{\Sigma} & \mathbf{0}_{m d^{\mathbf{R}} \times m' d^{\mathbf{R}'}} \\ \mathbf{0}_{m' d^{\mathbf{R}'}} \times m d^{\mathbf{R}} & \mathbf{I}_{m'} \otimes \boldsymbol{\Sigma}' \end{bmatrix} \right),$$

and

$$\mathbf{u} \equiv [\mathbf{u}_1^\top \dots \mathbf{u}_m^\top]^\top, \quad \mathbf{u}' \equiv [(\mathbf{u}'_1)^\top \dots (\mathbf{u}'_{m'})^\top]^\top.$$

The joining of the random effects vectors \mathbf{u} and \mathbf{u}' into one random vector from a single density function allows a form that is more amenable to EP. For likelihood inference with crossed random effects the integrals required grow with the number of random effects as well as the number of groups, i.e. when $d^{\mathbf{R}} = d^{\mathbf{R}'} = 1$ they are over the $m + m'$ -dimensional space. Although it is possible to reduce the integral to $m d^{\mathbf{R}}$ -dimensional inner integrals and outer integrals over $\mathbb{R}^{m' d^{\mathbf{R}'}}$ (alternatively $m' d^{\mathbf{R}'}$ -dimensional inner integrals and outer integrals over $\mathbb{R}^{m d^{\mathbf{R}}}$ is also possible), the dimension of the integral we wish to solve is still dependent on group size. Since the current form is more amenable to EP we do not conduct such a simplification here.

In this section we first explain the schematic of likelihood approximation using EP in Subsection [6.1.1](#), before exploring computation of point estimates and confidence intervals in Subsection [6.1.2](#). The results of our simulation studies are presented in Subsection [6.1.3](#).

6.1.1 Expectation propagation likelihood approximation

Although it seems that the extension of our methodology to crossed random effects is complicated, by expressing the likelihood as in equation [\(6.2\)](#) the updates for the projection onto the multivariate normal family from the general one level model apply to the crossed case. Specifically, EP for the crossed case centres around finding the

optimal natural parameters η_0 , η_1 and η_2 which minimise $\text{KL}(f_{\text{input}} \parallel f_{UN})$, where the approximating density f_{UN} is defined to be the unnormalised multivariate normal density in equation (3.4), and the target density $f_{\text{input}}(\mathbf{x})$ can be reparameterised to have the same form as the one level model

$$f_{\text{input}}(\mathbf{x}) = \Phi(c_0 + \mathbf{c}_1^\top \mathbf{x}) \exp\left(\left(\eta_1^{\text{input}}\right)^\top \mathbf{x} + \mathbf{x}^\top \mathbf{H}_2^{\text{input}} \mathbf{x}\right), \quad (6.3)$$

where η_1^{input} is a $d \times 1$ vector, $\mathbf{H}_2^{\text{input}}$ is a $d \times d$ matrix, $c_0 = (2y_{ii'j} - 1)\beta^\top \mathbf{x}_{ii'j}^{\mathbf{F}}$, $\mathbf{c}_1 = (2y_{ii'j} - 1)[(\mathbf{x}_{ii'j}^{\mathbf{R}})^\top (\mathbf{x}_{ii'j}^{\mathbf{R}'})^\top]^\top$ and $\mathbf{x} = [\mathbf{u}_i^\top (\mathbf{u}_{i'}^\top)]^\top$. It is easy to see that the moment matching problem required to find the optimal natural parameters to project onto the normalised multivariate normal family can be solved analogously to the one level probit model using Result 14, and the full projection onto the unnormalised multivariate normal family is given in Result 15.

6.1.1.1 Message passing formulation

The message passing scheme required for crossed random effects follows that of the one level model in Section 3.1.2, with minor ammendments. We explicitly show some of the updated messages here with the aim to clarify how the previous projection results are reused. Note the alternate expression of the prior

$$p(\mathbf{u}, \mathbf{u}'; \Sigma, \Sigma') \equiv \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \mathbf{u}_{i'} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}_{i'} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}_{i'} \end{array} \right]^\top \right) \end{array} \right]^\top \eta_{(\Sigma, \Sigma')} \right\},$$

where

$$\eta_{(\Sigma, \Sigma')} = \left[\begin{array}{c} -\frac{md^{\mathbf{R}}}{2} \log |2\pi \Sigma| - \frac{m'd^{\mathbf{R}'}}{2} \log |2\pi \Sigma'| \\ \mathbf{0}_{md^{\mathbf{R}}+m'd^{\mathbf{R}'}} \\ -\frac{1}{2} \mathbf{D}_{m'd^{\mathbf{R}'}+md^{\mathbf{R}}}^\top \text{vec} \left(\left[\begin{array}{cc} \mathbf{I}_m \otimes \Sigma & \mathbf{0}_{m'd^{\mathbf{R}} \times m'd^{\mathbf{R}'}} \\ \mathbf{0}_{m'd^{\mathbf{R}'} \times md^{\mathbf{R}}} & \mathbf{I}_{m'} \otimes \Sigma' \end{array} \right]^{-1} \right) \end{array} \right]. \quad (6.4)$$

The structure of equation (6.2) in factor graph form is shown in Figure 6.1, where the

circular stochastic node corresponds to the random vector $[\mathbf{u}^\top (\mathbf{u}')^\top]^\top$, the solid squares indicate the factor nodes and the dependencies of the factor nodes on the stochastic node are demonstrated through the edges.

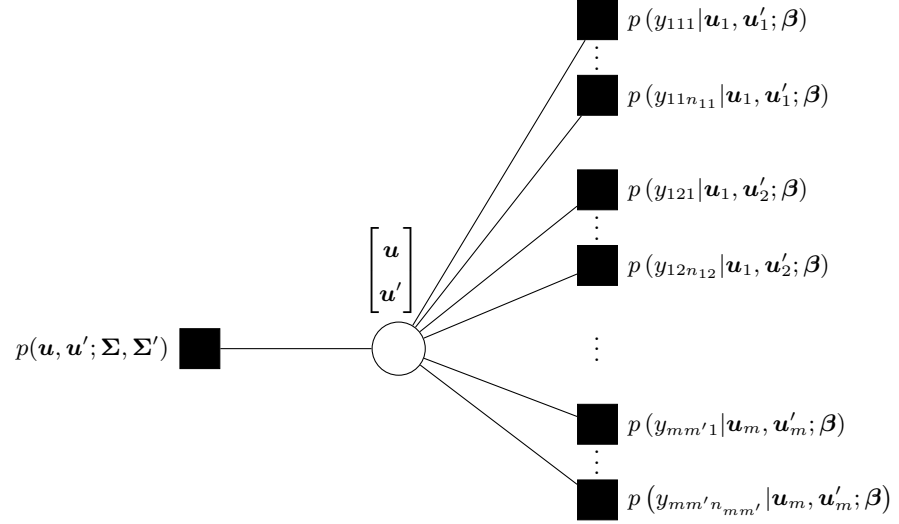


Figure 6.1: Factor graph representation of the product structure of the integrand in equation (6.2) for crossed random effects models.

Suppose that

$$p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) = \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \mathbf{u}'_{i'} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right]^\top \right) \end{array} \right]^\top \boldsymbol{\eta}_{ii'j} \right\}, \quad 1 \leq j \leq n_{ii'},$$

are initialised to be unnormalised multivariate normal density functions in $[\mathbf{u}_i^\top (\mathbf{u}'_{i'})^\top]^\top$. Then, for each $j = 1, \dots, n_{ii'}$, the $\boldsymbol{\eta}_{ii'j}$ update of the EP algorithm involves minimisation

of

$$\text{KL} \left(p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \prod_{j' \neq j}^{n_{ii'}} \tilde{p}(y_{ii'j'} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \times p(\mathbf{u}, \mathbf{u}'; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \left\| \prod_{(i,i') : n_{ii'} > 0} \prod_{j'=1}^{n_{ii'}} \tilde{p}(y_{ii'j'} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) p(\mathbf{u}, \mathbf{u}'; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \right. \right) \quad (6.5)$$

as a function of $(\mathbf{u}_i, \mathbf{u}'_{i'})$. As mentioned Result [15](#) can be used to update $\eta_{ii'j}$ in an iterative procedure until convergence.

The process of minimising the KL divergence in equation [\(6.5\)](#) can be expressed as a message from the factor $p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})$ to the stochastic node $[\mathbf{u}_i^\top (\mathbf{u}'_{i'})^\top]^\top$ as

$$\begin{aligned} & m_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \\ & \leftarrow \frac{\text{proj}_{UN} \left[m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})}(\mathbf{u}_i, \mathbf{u}'_{i'}) p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \right](\mathbf{u}_i, \mathbf{u}'_{i'})}{m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})}(\mathbf{u}_i, \mathbf{u}'_{i'})}, \end{aligned} \quad (6.6)$$

$$1 \leq i \leq m, \quad 1 \leq i' \leq m', \quad 1 \leq j \leq n_{ii'},$$

and the update of the message passed from $p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ to $(\mathbf{u}_i, \mathbf{u}'_{i'})$ is

$$\begin{aligned} & m_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \\ & \leftarrow \frac{\text{proj}_{UN} \left[m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}(\mathbf{u}_i, \mathbf{u}'_{i'}) p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \right](\mathbf{u}_i, \mathbf{u}'_{i'})}{m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}(\mathbf{u}_i, \mathbf{u}'_{i'})}. \end{aligned} \quad (6.7)$$

Similarly, the updates of stochastic node to factor messages are

$$\begin{aligned} & m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \\ & \leftarrow m_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \prod_{j' \neq j}^{n_{ii'}} m_{p(y_{ii'j'} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}), \end{aligned} \quad (6.8)$$

$$1 \leq i \leq m, \quad 1 \leq i' \leq m', \quad 1 \leq j \leq n_{ii'},$$

and

$$m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}(\mathbf{u}_i, \mathbf{u}'_{i'}) \leftarrow \prod_{j=1}^{n_{ii'}} m_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}). \quad (6.9)$$

We now express the key messages in their simplest natural parameter form. Recall

the natural parameter expression of the prior in equation (6.4) and that the unnormalised normal density constraint is enforced on equations (6.6) and (6.7). Then

$$m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}(\mathbf{u}_i, \mathbf{u}'_{i'}) \leftarrow \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \mathbf{u}'_{i'} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right]^\top \right) \end{array} \right]^\top \boldsymbol{\eta}_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')} \right\}. \quad (6.10)$$

By substituting the above forms into equation (6.7) it is easy to show the message $m_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \leftarrow p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ is constant throughout the message passing updates. As such, we now set

$$\boldsymbol{\eta}_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \leftarrow \boldsymbol{\eta}_{(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}. \quad (6.11)$$

For convenience, we denote the natural parameter vector

$$\boldsymbol{\eta}^\otimes \equiv \boldsymbol{\eta}_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})}.$$

Following the derivation of equation (6.10) it is easy to show

$$\begin{aligned} m_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})}(\mathbf{u}_i, \mathbf{u}'_{i'}) &= \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \mathbf{u}'_{i'} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right]^\top \right) \end{array} \right]^\top \boldsymbol{\eta}^\otimes \right\} \\ &= \exp(\eta_0^\otimes) \exp \left\{ \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right]^\top \boldsymbol{\eta}_1^\otimes + \left\{ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{array} \right]^\top \right) \right\}^\top \boldsymbol{\eta}_2^\otimes \right\}. \end{aligned}$$

Substituting this into equation (6.6) and following simplifications analogous to the

general one level case leads to

$$\begin{aligned}
& m_{p(y_{ij}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \\
& \leftarrow \text{proj}_{UN} \left[\Phi \left(c_{0_{ii'j}} + \mathbf{c}_{1_{ii'j}}^\top \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \right) \right. \\
& \quad \times \exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \boldsymbol{\eta}_1^\otimes + \left\{ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \right) \right\}^\top \boldsymbol{\eta}_2^\otimes \right\} \\
& \quad \times \left(\exp \left\{ \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \boldsymbol{\eta}_1^\otimes + \left\{ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \right) \right\}^\top \boldsymbol{\eta}_2^\otimes \right\} \right)^{-1},
\end{aligned}$$

where $c_{0_{ii'j}} = (2y_{ii'j} - 1)\boldsymbol{\beta}^\top \mathbf{x}_{ii'j}^{\mathbf{F}}$ and $\mathbf{c}_{1_{ii'j}} = (2y_{ii'j} - 1)[(\mathbf{x}_{ii'j}^{\mathbf{R}})^\top (\mathbf{x}_{ii'j}^{\mathbf{R}'})^\top]^\top$. Utilising Result [15](#),

$$\begin{aligned}
& m_{p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \\
& \leftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i \\ \mathbf{u}'_{i'} \\ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \right) \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right\},
\end{aligned}$$

where the linear and quadratic coefficient updates are

$$\left(\boldsymbol{\eta}_{p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_{1:2} \leftarrow K_{\text{probit}}(\boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ii'j}}, \mathbf{c}_{1_{ii'j}}) - \boldsymbol{\eta}_{1:2}^\otimes$$

and the constant coefficient update is

$$\begin{aligned}
& \left(\boldsymbol{\eta}_{p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_0 \\
& \leftarrow C_{\text{probit}} \left(\boldsymbol{\eta}_{1:2}^\otimes, \left(\boldsymbol{\eta}_{p(y_{ii'j}|\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_{1:2} + \boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ii'j}}, \mathbf{c}_{1_{ii'j}} \right).
\end{aligned}$$

Using the simplification of equation (6.6) and (6.7), equation (6.8) can be shown to be

$$m(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})(\mathbf{u}_i, \mathbf{u}'_{i'}) \leftarrow \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i \\ \mathbf{u}'_{i'} \\ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}^\top \right) \end{array} \right]^\top \boldsymbol{\eta}(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \right\},$$

where

$$\boldsymbol{\eta}(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \leftarrow \boldsymbol{\eta}_p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'}) + \sum_{j' \neq j} \boldsymbol{\eta}_p(y_{ii'j'} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'}).$$

Once convergence is reached, the EP approximation $\underline{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ is given by

$$\begin{aligned} \underline{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') &= \log \int_{\mathbb{R}^{md\mathbf{R}+m'd\mathbf{R}'}} \left(\prod_{(i,i')n_{ii'}>0} \prod_{j=1}^{n_{ii'}} m_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \right) \\ &\quad \times m_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) d \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix}. \end{aligned} \quad (6.12)$$

Using natural parameters the integral arising in equation (6.12) can be approximated by

$$\begin{aligned} &\int_{\mathbb{R}^{md\mathbf{R}+m'd\mathbf{R}'}} \left(\prod_{(i,i')n_{ii'}>0} \prod_{j=1}^{n_{ii'}} m_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) \right) \\ &\quad \times m_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}(\mathbf{u}_i, \mathbf{u}'_{i'}) d \begin{bmatrix} \mathbf{u}_i \\ \mathbf{u}'_{i'} \end{bmatrix} \\ &= (2\pi)^{(md\mathbf{R}+m'd\mathbf{R}')/2} \exp \left\{ \left(\boldsymbol{\eta}(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}') + \sum_{(i,i')n_{ii'}>0} \sum_{j=1}^{n_{ii'}} \boldsymbol{\eta}_p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'}) \right)_0 \right. \\ &\quad \left. + A \left(\left(\boldsymbol{\eta}(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}') + \sum_{(i,i')n_{ii'}>0} \sum_{j=1}^{n_{ii'}} \boldsymbol{\eta}_p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'}) \right)_{1:2} \right) \right\}. \end{aligned}$$

The full algorithm for the approximation of $\underline{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ using EP is provided in Algorithm

7

Algorithm 7 *Explicit form of algorithm used for the message passing approach to EP*

Inputs: $y_{ii'j}$, $\mathbf{x}_{ii'j}^{\mathbf{F}}$, $\mathbf{x}_{ii'j}^{\mathbf{R}}$, $\mathbf{x}_{ii'j}^{\mathbf{R}'}$, $1 \leq i \leq m$, $1 \leq i' \leq m'$, $1 \leq j \leq n_{ii'}$,

$\beta(d^{\mathbf{F}} \times 1)$, $\Sigma(d^{\mathbf{R}} \times d^{\mathbf{R}})$, $\Sigma'(d^{\mathbf{R}'} \times d^{\mathbf{R}'})$, are symmetric and positive definite.

Set constants: $c_{0_{ii'j}} \leftarrow (2y_{ii'j} - 1)\beta^{\top} \mathbf{x}_{ii'j}^{\mathbf{F}}$; $c_{1_{ii'j}} \leftarrow (2y_{ii'j} - 1)[(\mathbf{x}_{ii'j}^{\mathbf{R}})^{\top} (\mathbf{x}_{ii'j}^{\mathbf{R}'})^{\top}]^{\top}$,

$$1 \leq i \leq m, \quad 1 \leq i' \leq m', \quad 1 \leq j \leq n_{ii'};$$

$\eta_{p(\mathbf{u}_i, \mathbf{u}'_{i'}; \Sigma, \Sigma') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}$

$$\leftarrow \eta_{(\Sigma, \Sigma')} \equiv \left[\begin{array}{c} -\frac{m d^{\mathbf{R}}}{2} \log |2\pi \Sigma| - \frac{m' d^{\mathbf{R}'}}{2} \log |2\pi \Sigma'| \\ \mathbf{0}_{m d^{\mathbf{R}} + m' d^{\mathbf{R}'}} \\ -\frac{1}{2} \mathbf{D}^{\top} \left[\begin{array}{cc} \mathbf{I}_m \otimes \Sigma & \mathbf{0}_{m d^{\mathbf{R}} \times m' d^{\mathbf{R}'}} \\ \mathbf{0}_{m' d^{\mathbf{R}'} \times m d^{\mathbf{R}}} & \mathbf{I}_{m'} \otimes \Sigma' \end{array} \right]^{-1} \text{vec} \left(\begin{array}{c} \mathbf{0}_{m d^{\mathbf{R}} + m' d^{\mathbf{R}'}} \\ \mathbf{0}_{m' d^{\mathbf{R}' \times m d^{\mathbf{R}}} \\ \mathbf{I}_{m'} \otimes \Sigma' \end{array} \right) \end{array} \right],$$

$$1 \leq i \leq m, \quad 1 \leq i' \leq m'.$$

Initialise: $\eta_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}$, $1 \leq i \leq m$, $1 \leq i' \leq m'$, $1 \leq j \leq n_{ii'}$, using Laplace approximations.

Cycle:

$$\text{SUM} \left(\eta_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right) \leftarrow \sum_{(i, i') n_{ii'} > 0} \sum_{j=1}^{n_{ii'}} \eta_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}$$

For $(i, i') \in \{1, \dots, m\} \times \{1, \dots, m'\}$:

For $j = 1, \dots, n_{ii'}$:

$$\begin{aligned} \eta_{p(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta)} &\leftarrow \eta_{p(\mathbf{u}_i, \mathbf{u}'_{i'}, \Sigma, \Sigma') \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \\ &+ \text{SUM} \left(\eta_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right) - \eta_{p(y_{ij} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \\ &\left(\eta_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_{1:2} \\ &\leftarrow K_{\text{probit}} \left(\left(\eta_{\mathbf{u}_i \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta)} \right)_{1:2}; c_{0_{ii'j}}, c_{1_{ii'j}} \right) \\ &- \left(\eta_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \beta)} \right)_{1:2} \end{aligned}$$

until convergence of all natural parameters vectors.

Algorithm 8 Continuation of Algorithm [7](#) used for the message passing approach to EP

For $(i, i') \in \{1, \dots, m\} \times \{1, \dots, m'\}$:

For $j = 1, \dots, n_{ii'}$:

$$\begin{aligned} \left(\boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_0 &\leftarrow C_{\text{probit}} \left(\left(\boldsymbol{\eta}_{\mathbf{u}_i \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})} \right)_{1:2}, \right. \\ &\left. \left(\boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right)_{1:2} + \left(\boldsymbol{\eta}_{(\mathbf{u}_i, \mathbf{u}'_{i'}) \rightarrow p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta})} \right)_{1:2}; c_{0_{ii'j}}, c_{1_{ii'j}} \right). \end{aligned}$$

$$\text{SUM} \left(\boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right) \leftarrow \sum_{(i, i') n_{ii'} > 0} \sum_{j=1}^{n_{ii'}} \boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})}$$

Output: The full approximate log-likelihood is given by

$$\begin{aligned} \tilde{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}') &= \frac{1}{2} (m d^{\mathbf{R}} + m' d^{\mathbf{R}'}) \log(2\pi) + \left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}')} + \text{SUM} \left(\boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right) \right)_0 \\ &+ A \left\{ \left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}')} + \text{SUM} \left(\boldsymbol{\eta}_{p(y_{ii'j} | \mathbf{u}_i, \mathbf{u}'_{i'}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i, \mathbf{u}'_{i'})} \right) \right)_{1:2} \right\} \end{aligned}$$

where, $A(\boldsymbol{\eta})$ is defined as in equation [\(1.7\)](#) and $\boldsymbol{\eta}_{(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}')}$ follows from equation [\(2.21\)](#).

6.1.2 Computation of point estimates and confidence intervals

The numerical maximisation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}')$ is a simple extension of the one level case. As before, we must ensure the search for the maxima of $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ occur over the cone of symmetric positive definite matrices, which we accomplished with a re-parameterisation (Bateman & Pinheiro, 2000^{[55](#)}). For the general crossed random effects case the following procedure is recommended:

1. Convert $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ to the unconstrained space $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$ as follows:

(a) Obtain the spectral decomposition of $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$,

$$\boldsymbol{\Sigma} = \mathbf{u}_{\boldsymbol{\Sigma}} \text{diag}(\boldsymbol{\lambda}_{\boldsymbol{\Sigma}}) \mathbf{u}_{\boldsymbol{\Sigma}}^{\top} \quad \text{and} \quad \boldsymbol{\Sigma}' = \mathbf{u}'_{\boldsymbol{\Sigma}'} \text{diag}(\boldsymbol{\lambda}'_{\boldsymbol{\Sigma}'}) (\mathbf{u}'_{\boldsymbol{\Sigma}'})^{\top}.$$

(b) Used the spectral decomposition to obtain the matrix logarithm of $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$,

$$\log(\boldsymbol{\Sigma}) = \mathbf{u}_{\boldsymbol{\Sigma}} \text{diag}(\log(\boldsymbol{\lambda}_{\boldsymbol{\Sigma}})) \mathbf{u}_{\boldsymbol{\Sigma}}^{\top} \quad \text{and} \quad \log(\boldsymbol{\Sigma}') = \mathbf{u}'_{\boldsymbol{\Sigma}'} \text{diag}(\log(\boldsymbol{\lambda}'_{\boldsymbol{\Sigma}'})) (\mathbf{u}'_{\boldsymbol{\Sigma}'})^{\top}.$$

(c) Used $\log(\boldsymbol{\Sigma})$ and $\log(\boldsymbol{\Sigma}')$ to obtain $\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$ respectively,

$$\boldsymbol{\theta} \equiv \text{vech}\left(\frac{1}{2} \log(\boldsymbol{\Sigma})\right) \quad \text{and} \quad \boldsymbol{\theta}' \equiv \text{vech}\left(\frac{1}{2} \log(\boldsymbol{\Sigma}')\right).$$

2. Obtain the maximum likelihood estimate of $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}'})$ using a quasi-Newton optimisation method,

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}'}) = \text{argmax } \ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\theta}').$$

We conducted an initial search via Nelder-Mead with refinements by BFGS, which can be implemented via the “optim” R function.

3. Convert $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}'})$ to $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\boldsymbol{\omega}'})$ as follows:

(a) Obtain the spectral decomposition of $\text{vech}^{-1}(\hat{\boldsymbol{\theta}})$ and $\text{vech}^{-1}(\hat{\boldsymbol{\theta}'})$:

$$\text{vech}^{-1}(\hat{\boldsymbol{\theta}}) = \mathbf{u}_{\hat{\boldsymbol{\theta}}} \text{diag}(\boldsymbol{\lambda}_{\hat{\boldsymbol{\theta}}}) (\mathbf{u}_{\hat{\boldsymbol{\theta}}})^{\top} \quad \text{and} \quad \text{vech}^{-1}(\hat{\boldsymbol{\theta}'}) = \mathbf{u}'_{\hat{\boldsymbol{\theta}'}} \text{diag}(\boldsymbol{\lambda}'_{\hat{\boldsymbol{\theta}'}}) (\mathbf{u}'_{\hat{\boldsymbol{\theta}'}})^{\top}.$$

(b) Obtain $\hat{\boldsymbol{\Sigma}} = \mathbf{u}_{\hat{\boldsymbol{\theta}}} \text{diag}(\exp(2\boldsymbol{\lambda}_{\hat{\boldsymbol{\theta}}})) \mathbf{u}_{\hat{\boldsymbol{\theta}}}^{\top}$ and $\hat{\boldsymbol{\Sigma}}' = \mathbf{u}'_{\hat{\boldsymbol{\theta}'}} \text{diag}(\exp(2\boldsymbol{\lambda}'_{\hat{\boldsymbol{\theta}'}})) (\mathbf{u}'_{\hat{\boldsymbol{\theta}'}})^{\top}$.

(c) i. If $d^{\mathbf{R}} = 1$, then $\hat{\boldsymbol{\omega}} = \frac{1}{2} \log(\hat{\boldsymbol{\Sigma}})$. If $d^{\mathbf{R}'} = 1$, then $\hat{\boldsymbol{\omega}'} = \frac{1}{2} \log(\hat{\boldsymbol{\Sigma}'})$.
ii. If $d^{\mathbf{R}} > 1$, then

$$\hat{\boldsymbol{\omega}} = \left[\begin{array}{c} \frac{1}{2} \log(\text{diag}(\hat{\boldsymbol{\Sigma}})) \\ \tanh^{-1} \left\{ \frac{\text{vecbd}(\hat{\boldsymbol{\Sigma}})}{\sqrt{\text{vecbd}(\text{diag}(\hat{\boldsymbol{\Sigma}})\text{diag}(\hat{\boldsymbol{\Sigma}})^{\top})}} \right\} \end{array} \right].$$

If $d^{\mathbf{R}'} > 1$, then

$$\hat{\boldsymbol{\omega}'} = \left[\begin{array}{c} \frac{1}{2} \log(\text{diag}(\hat{\boldsymbol{\Sigma}'}) \\ \tanh^{-1} \left\{ \frac{\text{vecbd}(\hat{\boldsymbol{\Sigma}'})}{\sqrt{\text{vecbd}(\text{diag}(\hat{\boldsymbol{\Sigma}'})\text{diag}(\hat{\boldsymbol{\Sigma}'})^{\top})}} \right\} \end{array} \right].$$

4. Obtain the Hessian matrix $H\ell(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\boldsymbol{\omega}'})$ at the maximum $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\omega}}, \hat{\boldsymbol{\omega}'})$ using the quasi-Newton method BFGS, which as before can be implemented via `optim()`. The constraints on the parameters $(\boldsymbol{\beta}, \boldsymbol{\omega}, \boldsymbol{\omega}')$ means the Hessian should still be calculated on the $(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\theta}')$ space. The conversion can be obtained as follows:

(a) Form the symmetric matrices $\boldsymbol{\Sigma}$ and $\boldsymbol{\Sigma}'$ as follows:

i. If $d^{\mathbf{R}} = 1$, then $\boldsymbol{\Sigma} = \exp(2\boldsymbol{\omega})$. If $d^{\mathbf{R}'} = 1$ then $\boldsymbol{\Sigma}' = \exp(2\boldsymbol{\omega}')$.

ii. If $d^{\mathbf{R}} > 1$, then let $\boldsymbol{\omega}_1$ denote the first $d^{\mathbf{R}}$ entries of $\boldsymbol{\omega}$ and $\boldsymbol{\omega}_2$ denote the remaining $\frac{1}{2}d^{\mathbf{R}}(d^{\mathbf{R}} - 1)$ entries of $\boldsymbol{\omega}$. Set $\text{diag}(\boldsymbol{\Sigma}) = \exp(2\boldsymbol{\omega}_1)$. Obtain

the below-diagonal entries of Σ so that

$$\text{vecbd}(\Sigma) = \tanh(\omega_2) \odot \text{vecbd}(\exp(\omega_1) \exp(\omega_1)^\top)$$

holds. Obtain the above-diagonal entries of Σ such that symmetry of Σ is enforced. If $d^{\mathbf{R}'} > 1$, then let ω'_1 denote the first $d^{\mathbf{R}'}$ entries of ω' and ω'_2 denote the remaining $\frac{1}{2}d^{\mathbf{R}'}(d^{\mathbf{R}'} - 1)$ entries of ω' . Set $\text{diag}(\Sigma') = \exp(2\omega'_1)$. Obtain the below-diagonal entries of Σ' so that

$$\text{vecbd}(\Sigma') = \tanh(\omega'_2) \odot \text{vecbd}(\exp(\omega'_1) \exp(\omega'_1)^\top)$$

holds. Obtain the above-diagonal entries of Σ' such that symmetry of Σ' is enforced.

iii. Obtain the spectral decomposition:

$$\Sigma = \mathbf{u}_\Sigma \text{diag}(\lambda_\Sigma) \mathbf{u}_\Sigma^\top \quad \text{and} \quad \Sigma' = \mathbf{u}'_{\Sigma'} \text{diag}(\lambda'_{\Sigma'}) (\mathbf{u}'_{\Sigma'})^\top.$$

iv. Obtain

$$\theta = \text{vech} \left\{ \frac{1}{2} \mathbf{u}_\Sigma \text{diag}(\log(\lambda_\Sigma)) \mathbf{u}_\Sigma^\top \right\}$$

and

$$\theta' = \text{vech} \left\{ \frac{1}{2} \mathbf{u}'_{\Sigma'} \text{diag}(\log(\lambda'_{\Sigma'})) (\mathbf{u}'_{\Sigma'})^\top \right\}.$$

5. Form $100(1 - \alpha)\%$ confidence intervals for the entries of (β, ω, ω') using

$$\begin{bmatrix} \hat{\beta} \\ \hat{\omega} \\ \hat{\omega}' \end{bmatrix} \pm \Phi^{-1} \left(1 - \frac{1}{2} \alpha \right) \sqrt{-\text{diag} \left\{ (\mathbf{H} \ell(\hat{\beta}, \hat{\omega}, \hat{\omega}'))^{-1} \right\}}.$$

6. Back transform the confidence interval limits for the ω component, to correspond to the standard deviation and correlation parameters as follows:

$$\begin{bmatrix} \sqrt{\text{diag}(\Sigma)} \\ \text{vecbd}(\Sigma) / \sqrt{\text{vecbd}(\text{diag}(\Sigma) \text{diag}(\Sigma)^\top)} \end{bmatrix}$$

and

$$\begin{bmatrix} \sqrt{\text{diag}(\Sigma')} \\ \text{vecbd}(\Sigma') / \sqrt{\text{vecbd}(\text{diag}(\Sigma') \text{diag}(\Sigma')^\top)} \end{bmatrix}.$$

6.1.3 Simulation study

Three simulation studies were conducted, each for 1000 replicates on University of Technology Sydney Interactive High Performance Computing facility Jupiter node with eight 3.6 gigahertz processors and 32 gigabytes random access memory.

6.1.3.1 Comparison with MCMC and Laplace approximation maximum likelihood for crossed random effects

The observations were simulated 1000 times according to equation (6.1) with true parameter values

$$\boldsymbol{\beta}_{\text{true}} = [-0.58, 1.07]^\top, \quad \boldsymbol{\Sigma}_{\text{true}} = \sigma_{\text{true}}^2 = 0.32, \quad \text{and} \quad \boldsymbol{\Sigma}'_{\text{true}} = (\sigma_{\text{true}}^2)' = 0.47. \quad (6.13)$$

The number of groups in the data were $m = 10$, $m' = 6$ with the number of measurements in each group fixed at $n_{ii'} = 3$ for all (i, i') . The $\mathbf{x}_{ii'j}^{\mathbf{F}}$, $\mathbf{x}_{ii'j}^{\mathbf{R}}$ and $\mathbf{x}_{ii'j}^{\mathbf{R}'}$ vectors were of the form

$$\mathbf{x}_{ii'j}^{\mathbf{F}} = [1, x_{ii'j}]^\top, \quad \mathbf{x}_{ii'j}^{\mathbf{R}} = 1 \quad \text{and} \quad \mathbf{x}_{ii'j}^{\mathbf{R}'} = 1$$

where $x_{ii'j}$ was generated independently from a uniform distribution on the unit interval. The tolerance of error values for the EP algorithm was set to 10^{-5} and a maximum of 100 Nelder-Mead search iterations was used. The model described was fit using each of the following methods:

1. Laplace approximation implemented via the function `glmer()` in the R package “lme4” (Bates, et al., 2018^[5]).
2. EP as described in this Section, initialised using Laplace approximations from the function `glmer()` in the R package “lme4” (Bates, et al., 2018^[5]).
3. Markov Chain Monte Carlo with flat priors.

We compare point estimates and confidence intervals produced by each of the three methods in Figure 6.2. Our studies show that the quality of results produced by each method varied across parameters. MCMC provided the highest coverage (98.6%) for the fixed intercept, while EP and Laplace approximations provided estimates which were over confident (92.7% and 93.2% respectively). For the fixed slope all methods provided approximately 95% coverage, with Laplace approximations providing the highest empirical coverage (95.1%). No method provided good results for variance

parameter estimates. Laplace approximations had the highest empirical coverage for the variance parameter of the first level (69.5% compared to 63% from EP), whilst EP had the highest empirical coverage for the variance parameter of the crossed level (89.2% compared to 85.9% from Laplace approximations). MCMC provided the worst coverage for the first and crossed level of (56.5% and 69.2% respectively). Overall it seems that Laplace approximations perform marginally better than EP for this model and data senario.

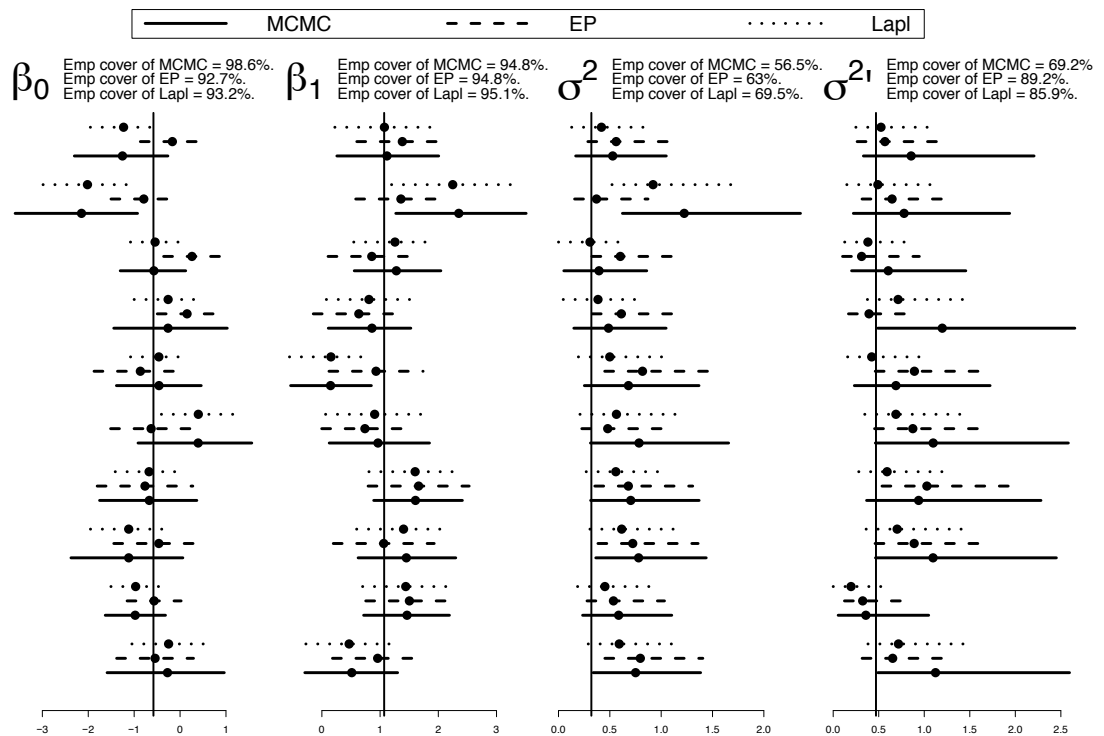


Figure 6.2: Comparison of point estimates and 95% confidence intervals for the simulation study with true parameter values given by equation (6.13). We display 10 randomly chosen replications of the simulation study described. The vertical lines indicate true parameter values and the percentages displayed at the top of each panel are empirical coverages over all 1000 replications for each method involved in the comparison.

6.2 The general probit two level mixed model

We now extend our methodology for the probit case to deal with two level GLMMs, utilising key results from Chapter 3 as in the crossed case. Our aim is to find an approximation to the maximum likelihood estimates of the parameters β , Σ^{L1} and Σ^{L2}

with 95% confidence intervals.

The probit model with two levels of nesting is specified in equation (6.14),

$$\begin{aligned}
y_{ijk} | \mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left\{ \Phi \left(\boldsymbol{\beta}^\top \mathbf{x}_{ijk}^{\mathbf{F}} + (\mathbf{u}_i^{\mathbf{L1}})^\top \mathbf{x}_{ijk}^{\mathbf{R1}} + (\mathbf{u}_{ij}^{\mathbf{L2}})^\top \mathbf{x}_{ijk}^{\mathbf{R2}} \right) \right\}, \\
\mathbf{u}_i^{\mathbf{L1}} &\stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R1}}}, \boldsymbol{\Sigma}^{\mathbf{L1}}) \quad \text{independently of} \quad \mathbf{u}_{ij}^{\mathbf{L2}} \stackrel{\text{ind.}}{\sim} \mathbf{N}(\mathbf{0}_{d^{\mathbf{R2}}}, \boldsymbol{\Sigma}^{\mathbf{L2}}), \\
1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij}, & \tag{6.14}
\end{aligned}$$

where the response $y_{ijk} \in \{0, 1\}$ and predictor vectors $\mathbf{x}_{ijk}^{\mathbf{F}}$, $\mathbf{x}_{ijk}^{\mathbf{R1}}$ and $\mathbf{x}_{ijk}^{\mathbf{R2}}$ correspond to the k th set of measurements within the j th inner group within the i th outer group. The number of outer groups is m and the number of inner groups in the i th outer group is n_i . The sample size of the j th group in the i th outer group is o_{ij} . $\mathbf{x}_{ijk}^{\mathbf{R1}}$ and $\mathbf{x}_{ijk}^{\mathbf{R2}}$ have respective dimensions $d^{\mathbf{R1}} \times 1$ and $d^{\mathbf{R2}} \times 1$.

The form of the log-likelihood is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}) \equiv \sum_{i=1}^m \ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}),$$

where

$$\begin{aligned}
\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}) &= \log \int_{\mathbb{R}^{d^{\mathbf{R1}}+d^{\mathbf{R2}}}} \prod_{j=1}^{n_i} \prod_{k=1}^{o_{ij}} p(y_{ijk} | \mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}}; \boldsymbol{\beta}) \\
&\quad \times p(\mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}}; \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}) d \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}, \tag{6.15}
\end{aligned}$$

$$p(y_{ijk} | \mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}}; \boldsymbol{\beta}) = \Phi \left\{ (2y_{ijk} - 1) \left(\boldsymbol{\beta}^\top \mathbf{x}_{ijk}^{\mathbf{F}} + (\mathbf{u}_i^{\mathbf{L1}})^\top \mathbf{x}_{ijk}^{\mathbf{R1}} + (\mathbf{u}_{ij}^{\mathbf{L2}})^\top \mathbf{x}_{ijk}^{\mathbf{R2}} \right) \right\}$$

and

$$\begin{aligned}
p(\mathbf{u}_i^{\mathbf{L1}}, \mathbf{u}_{ij}^{\mathbf{L2}}; \boldsymbol{\Sigma}^{\mathbf{L1}}, \boldsymbol{\Sigma}^{\mathbf{L2}}) &\equiv |2\pi \boldsymbol{\Sigma}^{\mathbf{L1}}|^{-1/2} |2\pi \boldsymbol{\Sigma}^{\mathbf{L2}}|^{-1/2} \\
&\quad \times \exp \left\{ -\frac{1}{2} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Sigma}^{\mathbf{L1}} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}^{\mathbf{L2}} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L1}} \\ \mathbf{u}_{ij}^{\mathbf{L2}} \end{bmatrix} \right\}.
\end{aligned}$$

Unlike the crossed random effects case, the integrals required are not dependent on the number of groups. Subsection 6.2.1 explains the schematic of likelihood approximation for the two level probit model using EP. It is easy to see how computation of point

estimates and confidence intervals for $(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\text{L1}}, \boldsymbol{\Sigma}^{\text{L2}})$ using `optim()` in R follows the same process as the crossed model in Subsection 6.1.2, and as such we do not repeat it. Due to time constraints we did not conduct a simulation study and as such no results are presented.

6.2.1 Expectation propagation likelihood approximation

As in the crossed random effects case, the results facilitating the projection of the target density onto the unnormalised multivariate normal family for the general one level probit model can be reused for the two level probit model. For the two level model EP involves finding the optimal natural parameters $\eta_0, \boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ which minimise $\text{KL}(f_{\text{input}} \parallel f_{\text{UN}})$, where the unnormalised multivariate normal density f_{UN} is defined by equation 3.4 and the target density $f_{\text{input}}(\boldsymbol{x})$ can be written as

$$f_{\text{input}}(\boldsymbol{x}) = \Phi(c_0 + \boldsymbol{c}_1^\top \boldsymbol{x}) \exp\left(\left(\boldsymbol{\eta}_1^{\text{input}}\right)^\top \boldsymbol{x} + \boldsymbol{x}^\top \boldsymbol{H}_2^{\text{input}} \boldsymbol{x}\right), \quad (6.16)$$

where $\boldsymbol{\eta}_1^{\text{input}}$ is a $d \times 1$ vector, $\boldsymbol{H}_2^{\text{input}}$ is a $d \times d$ matrix, $c_0 = (2y_{ijk} - 1)\boldsymbol{\beta}^\top \boldsymbol{x}_{ijk}^{\text{F}}$, $\boldsymbol{c}_1 = (2y_{ijk} - 1)[(\boldsymbol{x}_{ijk}^{\text{R1}})^\top (\boldsymbol{x}_{ijk}^{\text{R2}})^\top]^\top$ and $\boldsymbol{x} = [(\boldsymbol{u}_i^{\text{L1}})^\top (\boldsymbol{u}_{ij}^{\text{L2}})^\top]^\top$. It is easy to see we require $\boldsymbol{\eta}^*$ to solve equation 3.7, where the optimal natural parameters $\boldsymbol{\eta}_1^*$ and $\boldsymbol{\eta}_2^*$ to project onto the multivariate normal family are given by Result 14 and the projection on the unnormalised multivariate normal density is given by Result 15.

6.2.1.1 Message passing formulation

The message passing scheme required for the two level model follows that of Section 3.1.2, with minor amendments. For clarity we explicitly show it here. Note the alternate matrix expression of the prior

$$p(\boldsymbol{u}_i^{\text{L1}}, \boldsymbol{u}_{ij}^{\text{L2}}, \boldsymbol{\Sigma}^{\text{L1}}, \boldsymbol{\Sigma}^{\text{L2}}) \equiv \exp \left\{ \left[\begin{array}{c} 1 \\ \boldsymbol{u}_i^{\text{L1}} \\ \boldsymbol{u}_{ij}^{\text{L2}} \\ \text{vech} \left(\begin{bmatrix} \boldsymbol{u}_i^{\text{L1}} \\ \boldsymbol{u}_{ij}^{\text{L2}} \end{bmatrix} \begin{bmatrix} \boldsymbol{u}_i^{\text{L1}} \\ \boldsymbol{u}_{ij}^{\text{L2}} \end{bmatrix}^\top \right) \end{array} \right]^\top \boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\text{L1}}, \boldsymbol{\Sigma}^{\text{L2}})} \right\},$$

where

$$\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} \equiv \begin{bmatrix} -\frac{d^{\mathbf{R}1}}{2} \log |2\pi \boldsymbol{\Sigma}^{\mathbf{L}1}| - \frac{d^{\mathbf{R}2}}{2} \log |2\pi \boldsymbol{\Sigma}^{\mathbf{L}2}| \\ \mathbf{0}_{d^{\mathbf{R}1}+d^{\mathbf{R}2}} \\ -\frac{1}{2} \mathbf{D}_{d^{\mathbf{R}1}+d^{\mathbf{R}2}}^\top \text{vec} \left(\begin{bmatrix} \boldsymbol{\Sigma}^{\mathbf{L}1} & \mathbf{0}_{d^{\mathbf{R}1} \times d^{\mathbf{R}2}} \\ \mathbf{0}_{d^{\mathbf{R}2} \times d^{\mathbf{R}1}} & \boldsymbol{\Sigma}^{\mathbf{L}2} \end{bmatrix}^{-1} \right) \end{bmatrix}. \quad (6.17)$$

The structure of equation (6.2) in factor graph form is shown in Figure 6.3, where the circular stochastic node corresponds to the random vector $(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})$, the solid squares indicate the factor nodes and the dependencies of the factor nodes on the stochastic node are demonstrated through the edges. Suppose that

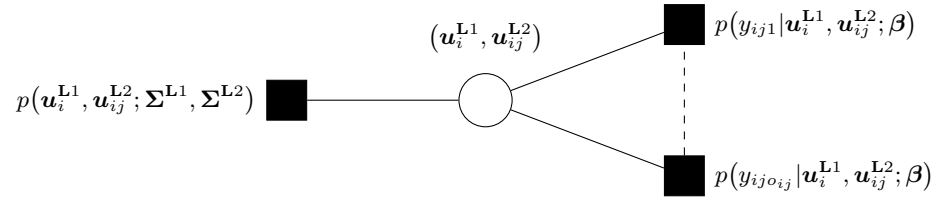


Figure 6.3: Factor graph representation of the product structure of the integrand in equation (6.15) for random effects models with two levels of nesting.

$$p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) = \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \\ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \right) \end{bmatrix}^\top \boldsymbol{\eta}_{ijk} \right\}, \quad 1 \leq k \leq o_{ij},$$

are initialised to be unnormalised multivariate normal density functions in $(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})$. Then, for each $k = 1, \dots, o_{ij}$, the $\boldsymbol{\eta}_{ijk}$ update involves minimisation of

$$\begin{aligned} & \text{KL} \left(p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \prod_{j=1}^{n_i} \prod_{k' \neq k}^{o_{ij}} p(y_{ijk'} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \right. \\ & \left. \times p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \right\| \prod_{j=1}^{n_i} \prod_{k'=1}^{o_{ij}} p(y_{ijk'} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \end{aligned} \quad (6.18)$$

as a function of $(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})$. We can achieve this by using Result [15](#) to update η_{ijk} in an iterative procedure until it converges.

The process of minimising the KL divergence in equation [\(6.18\)](#) can be expressed as a message from the factor $p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)$ to the stochastic node $(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})$ as

$$m_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \leftarrow \frac{\text{proj}_{\mathcal{UN}} \left[m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \right] (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}{m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}, \quad (6.19)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij},$$

and the update of the message passed from $p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2})$ to $(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})$ is

$$m_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \leftarrow \frac{\text{proj}_{\mathcal{UN}} \left[m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2}) \right] (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}{m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}. \quad (6.20)$$

The updates of stochastic node to factor messages are

$$m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) = m_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \prod_{j=1}^{n_i} \prod_{k' \neq k}^{o_{ij}} m_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}), \quad (6.21)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ijk},$$

and

$$m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) = \prod_{j=1}^{n_i} \prod_{k'=1}^{o_{ij}} m_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}). \quad (6.22)$$

Recall the natural parameter expression of the prior in equation [\(6.17\)](#) and that the unnormalised normal density constraint is enforced on equations [\(6.19\)](#) and [\(6.20\)](#).

Then

$$m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) = \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right]^{\top} \right) \end{array} \right]^{\top} \boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} \right\}. \quad (6.23)$$

By substituting the above forms into equation (6.20) it is easy to show the message $m_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \leftarrow p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})$ is constant throughout the message passing updates. As such, we now set

$$\boldsymbol{\eta}_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \leftarrow \boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})}. \quad (6.24)$$

For convenience, we denote the natural parameter vector

$$\boldsymbol{\eta}^{\otimes} \equiv \boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})}.$$

Following the derivation of equation (6.23) it is easy to show

$$\begin{aligned} m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) &= \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \\ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right]^{\top} \right) \end{array} \right]^{\top} \boldsymbol{\eta}^{\otimes} \right\} \\ &= \exp(\eta_0^{\otimes}) \exp \left\{ \left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right]^{\top} \boldsymbol{\eta}_1^{\otimes} + \left\{ \text{vech} \left(\left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{array} \right]^{\top} \right) \right\}^{\top} \boldsymbol{\eta}_2^{\otimes} \right\}. \end{aligned}$$

Substituting this into equation (6.19) and following simplifications analogous to the

general one level case leads to

$$\begin{aligned}
& m_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \\
& \leftarrow \text{proj}_{UN} \left[\Phi \left(c_{0_{ijk}} + \mathbf{c}_{1_{ijk}}^\top \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \right) \right. \\
& \quad \times \exp \left\{ \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \boldsymbol{\eta}_1^\otimes + \left\{ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \right) \right\}^\top \boldsymbol{\eta}_2^\otimes \right\} \\
& \quad \times \left(\exp \left\{ \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \boldsymbol{\eta}_1^\otimes + \left\{ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \right) \right\}^\top \boldsymbol{\eta}_2^\otimes \right\} \right)^{-1},
\end{aligned}$$

where $c_{0_{ijk}} = (2y_{ijk} - 1)\boldsymbol{\beta}^\top \mathbf{x}_{ijk}^{\mathbf{F}}$ and $\mathbf{c}_{1_{ijk}} = (2y_{ijk} - 1)[(\mathbf{x}_{ijk}^{\mathbf{R}1})^\top (\mathbf{x}_{ijk}^{\mathbf{R}2})^\top]^\top$. Utilising Result [15](#),

$$\begin{aligned}
& m_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \\
& \leftarrow \exp \left\{ \begin{bmatrix} 1 \\ \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \\ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \right) \end{bmatrix}^\top \boldsymbol{\eta}_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right\},
\end{aligned}$$

where the linear and quadratic coefficient updates are

$$\left(\boldsymbol{\eta}_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_{1:2} \leftarrow K_{\text{probit}}(\boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ijk}}, \mathbf{c}_{1_{ijk}}) - \boldsymbol{\eta}_{1:2}^\otimes$$

and the constant coefficient update is

$$\begin{aligned}
\left(\boldsymbol{\eta}_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_0 & \leftarrow C_{\text{probit}} \left(\boldsymbol{\eta}_{1:2}^\otimes, \left(\boldsymbol{\eta}_{p(y_{ijk}|\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_{1:2} \right. \\
& \quad \left. + \boldsymbol{\eta}_{1:2}^\otimes; c_{0_{ijk}}, \mathbf{c}_{1_{ijk}} \right).
\end{aligned}$$

Using the simplification of equation [\(6.19\)](#) and [\(6.20\)](#), equation [\(6.21\)](#) can be shown to

be

$$m_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \leftarrow \exp \left\{ \left[\begin{array}{c} 1 \\ \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \\ \text{vech} \left(\begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}^\top \right) \end{array} \right]^\top \boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})} \right\},$$

where

$$\begin{aligned} \boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \\ &+ \sum_{j' \neq j} \boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij'}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}. \end{aligned}$$

Once convergence is reached, the EP approximation of each $\ell_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})$ is given by

$$\begin{aligned} \tilde{\ell}_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) &= \log \int_{\mathbb{R}^{d\mathbf{R}1+d\mathbf{R}2}} \prod_{j=1}^{n_i} \prod_{k=1}^{o_{ij}} \left(m_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \right) \\ &\times m_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) d \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix}. \end{aligned} \quad (6.25)$$

Using natural parameters the integral arising in equation (6.25) can be approximated by

$$\begin{aligned} &\int_{\mathbb{R}^{d\mathbf{R}1+d\mathbf{R}2}} \prod_{j=1}^{n_i} \prod_{k=1}^{o_{ij}} \left(m_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \right) \\ &\times m_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) d \begin{bmatrix} \mathbf{u}_i^{\mathbf{L}1} \\ \mathbf{u}_{ij}^{\mathbf{L}2} \end{bmatrix} \\ &= (2\pi)^{(d\mathbf{R}1+d\mathbf{R}2)/2} \exp \left\{ \left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} + \sum_{j=1}^{n_i} \sum_{k=1}^{o_{ij}} \boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_0 \right. \\ &\left. + A \left(\left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} + \sum_{j=1}^{n_i} \sum_{k=1}^{o_{ij}} \boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_{1:2} \right) \right\}. \end{aligned}$$

The full algorithm for the approximation of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})$ using EP is provided in

Algorithm [9](#)

Algorithm 9 *Explicit form of algorithm used for the message passing approach to EP*

Inputs: y_{ijk} , $\mathbf{x}_{ijk}^{\mathbf{F}}$, $\mathbf{x}_{ijk}^{\mathbf{R}1}$, $\mathbf{x}_{ijk}^{\mathbf{R}2}$, $1 \leq i \leq m$, $1 \leq j \leq n_i$, $1 \leq k \leq o_{ij}$,

$\beta(d^{\mathbf{F}} \times 1)$, $\Sigma^{\mathbf{L}1}(d^{\mathbf{R}1} \times d^{\mathbf{R}1})$, $\Sigma^{\mathbf{L}2}(d^{\mathbf{R}2} \times d^{\mathbf{R}2})$, are symmetric and positive definite.

Set constants: $c_{0_{ijk}} \leftarrow (2y_{ijk} - 1)\beta^{\top} \mathbf{x}_{ijk}^{\mathbf{F}}$; $\mathbf{c}_{1_{ijk}} \leftarrow (2y_{ijk} - 1)[(\mathbf{x}_{ijk}^{\mathbf{R}1})^{\top} (\mathbf{x}_{ijk}^{\mathbf{R}2})^{\top}]^{\top}$,

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij};$$

$\eta_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma, \Sigma') \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}$

$$\leftarrow \eta_{(\Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2})} \equiv \begin{bmatrix} -\frac{d^{\mathbf{R}1}}{2} \log|2\pi \Sigma^{\mathbf{L}1}| - \frac{d^{\mathbf{R}2}}{2} \log|2\pi \Sigma^{\mathbf{L}2}| \\ \mathbf{0}_{d^{\mathbf{R}1}+d^{\mathbf{R}2}} \\ -\frac{1}{2} \mathbf{D}_{d^{\mathbf{R}1}+d^{\mathbf{R}2}}^{\top} \text{vec} \left(\begin{bmatrix} \Sigma^{\mathbf{L}1} & \mathbf{0}_{d^{\mathbf{R}1} \times d^{\mathbf{R}2}} \\ \mathbf{0}_{d^{\mathbf{R}2} \times d^{\mathbf{R}1}} & \Sigma^{\mathbf{L}2} \end{bmatrix}^{-1} \right) \end{bmatrix},$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij}.$$

For $i = 1, \dots, m$:

Initialise: $\eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}$, $1 \leq i \leq m$, $1 \leq j \leq n_i$, $1 \leq k \leq o_{ij}$, using Laplace approximations.

Cycle:

$$\text{SUM} \left(\eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right) \leftarrow \sum_{j=1}^{n_i} \sum_{k=1}^{o_{ij}} \eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}$$

For $j = 1, \dots, n_i$:

For $k = 1, \dots, o_{ij}$:

$$\begin{aligned} \eta_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)} &\leftarrow \eta_{p(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \Sigma^{\mathbf{L}1}, \Sigma^{\mathbf{L}2}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \\ &+ \text{SUM} \left(\eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right) - \eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \\ &\left(\eta_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_{1:2} \\ &\leftarrow K_{\text{probit}} \left(\left(\eta_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)} \right)_{1:2}; c_{0_{ijk}}, \mathbf{c}_{1_{ijk}} \right) \\ &- \left(\eta_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \beta)} \right)_{1:2} \end{aligned}$$

until convergence of all natural parameters vectors.

Algorithm 10 Continuation of Algorithm 9 used for the message passing approach to EP

For $j = 1, \dots, n_i$:

For $k = 1, \dots, o_{ij}$:

$$\left(\boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_0 \leftarrow C_{\text{probit}} \left(\left(\boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})} \right)_{1:2}, \right. \\ \left. \left(\boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right)_{1:2} + \left(\boldsymbol{\eta}_{(\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}) \rightarrow p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta})} \right)_{1:2}; \right. \\ \left. c_{0_{ijk}}, c_{1_{ijk}} \right).$$

$$\text{SUM} \left(\boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right) \leftarrow \sum_{j=1}^{n_i} \sum_{k=1}^{o_{ij}} \boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})}$$

Output: The full approximate log-likelihood is given by

$$\tilde{\ell}(\boldsymbol{\beta}, \boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2}) = \frac{m(d^{\mathbf{R}1} + d^{\mathbf{R}2})}{2} \log(2\pi) \\ + \sum_{i=1}^m \left\{ \left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} + \text{SUM} \left(\boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right) \right)_0 \right. \\ \left. + A \left\{ \left(\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})} + \text{SUM} \left(\boldsymbol{\eta}_{p(y_{ijk} | \mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2}; \boldsymbol{\beta}) \rightarrow (\mathbf{u}_i^{\mathbf{L}1}, \mathbf{u}_{ij}^{\mathbf{L}2})} \right) \right)_{1:2} \right\} \right\},$$

where $A(\boldsymbol{\eta})$ is defined as in equation (1.7) and $\boldsymbol{\eta}_{(\boldsymbol{\Sigma}^{\mathbf{L}1}, \boldsymbol{\Sigma}^{\mathbf{L}2})}$ follows from equation (2.21).

Chapter 7

Applications of expectation propagation for one level probit mixed models

Regardless of the amount of care taken during data collection, real data is subject to the laws of the environment it is collected in; that is, it is often full of subtle yet difficult features and traits, providing unforeseen challenges to analysts when fitting models. The facets of real data are not easily replicated in simulations due to the stochastic environment in which they are created. As such, each real dataset provides a unique opportunity for rigorous testing of statistical methodology. Not only does this help to expose unknown weaknesses which can be either acknowledged or further developed, but can also help highlight strengths. In Section [7.1](#) we introduce the R package created for implementing the methodology in the probit case as explained in Chapter [3](#). Section [7.2](#) and [7.3](#) this cover applications of our methodology on two real datasets, each with unique challenges. We fit models for each dataset and compare the fits to current methodology, and in the case of the second dataset contribute some useful findings.

7.1 R package “glmmEP”

To facilitate implementation of the methodology presented in this thesis we developed the R package “glmmEP” (Wand & Yu, 2020^{[69](#)}). Due to computational issues with the integrals required to solve non-probit models and time constraints, the package is restricted to one level probit models. Although the package uses an R interface, all computations relating to EP are conducted in Fortran77 to reduce computing

time. The R function `glmEP()` is used to fit a probit GLMM using our methodology, where the user must supply a vector of binary responses, a matrix of fixed effects, a matrix of random effects, and a vector with identification numbers for each group. Calling `summary.glmEP()` returns the maximum likelihood estimates and corresponding confidence intervals for the fixed effects parameters. Confidence levels are set using the `confLev` argument in the function `glmEP.control()`. Unlike other packages we also provide confidence intervals for random effects parameters. A vignette with more details is available by calling `glmEPvignette()`.

7.2 Modelling immunisation of Guatemalan children

In 1987, the National Survey of Maternal and Child Health was conducted in Guatemala on a sample of 5160 women age 15 to 44. The study utilised a questionnaire about prenatal care and immunisation status of children born during the study period. Pebley, Goldman & Rodríguez (1996)⁵⁴ provide a thorough explanation of the 1987 study dataset and conduct an analysis.

We used the dataset `guImmun` from the R package “`mlmRev`” (Bates, Maechler & Bolker, 2019)⁴ which is a smaller version of the full dataset of the 1987 National Survey of Maternal and Child Health, with only 2519 observations and 13 variables. Of the 13 we selected 8, which are summarised in Table (7.1). This dataset has a low number of observations per grouping variable and as such is a difficult model to estimate. Our aim is to compare the fits of current GLMM methodology to that provided by our methodology.

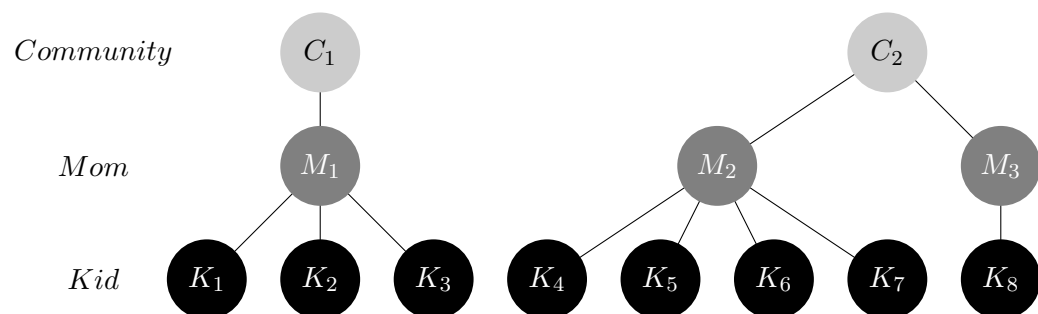


Figure 7.1: A plot showing the two level multiple membership structure of the data, specifically kids nested under mothers nested under community. The first row shows each community, the second shows each mother in each community and the third row shows the kids of each mother in each community.

Table 7.1: *A list of the variables in the `guImmun` dataset used for model fitting.*

Variable name	Description
<code>immun</code>	A two-level factor whether or not a child has their complete set of immunisations.
<code>pcInd81</code>	A continuous variable of indigenous population percentage in the community child lived in during 1981 census.
<code>kid2p</code>	A two-level factor specifying if the child was two years or older.
<code>momEd</code>	A three-level factor with the mother's level of school education. The levels are not finished primary school, finished primary school, and finished secondary school.
<code>husEd</code>	A four-level factor variable with the husband's level of school education. The levels are not finished primary school, finished primary school, finished secondary school, and unknown.
<code>momWork</code>	A two-level factor variable whether the child's mother had ever worked outside the home.
<code>rural</code>	A two-level factor variable indicating whether or not the child's location is rural or urban.
<code>mom</code>	A multilevel factor variable that codes the children's mothers.

The random intercepts and slopes probit mixed model we fit follows in equation (7.1),

$$\begin{aligned} \mathcal{I}(\text{immun}_{ij} = Y) | u_{0i}, u_{1i} \stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left\{ \Phi \left(\beta_0 + u_{0i} + (\beta_1 + u_{1i}) \text{pcInd81}_{ij} \right. \right. \\ \left. \left. + \beta_2 \mathcal{I}(\text{kid2p}_{ij} = Y) + \beta_3 \mathcal{I}(\text{momEd}_{ij} = S) + \beta_4 \mathcal{I}(\text{husEd}_{ij} = S) \right. \right. \\ \left. \left. + \beta_5 \mathcal{I}(\text{momWork}_{ij} = Y) + \beta_6 \mathcal{I}(\text{rural}_{ij} = Y) \right) \right\}, \end{aligned} \quad (7.1)$$

where $\mathcal{I}(\mathcal{P}) = 1$ when \mathcal{P} is true and 0 otherwise, and immun_{ij} is the immunity value for the j th child of the i th mother for $1 \leq i \leq 1596$, $1 \leq j \leq n_i$ and $n_i \in \{1, 2, 3\}$. Other variables in the model are defined analogously. We assume our bivariate random effects vectors satisfy

$$\begin{bmatrix} u_{0i} \\ u_{1i} \end{bmatrix} \stackrel{\text{ind.}}{\sim} \mathbf{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right). \quad (7.2)$$

Note that although the data does have an additional level of nesting we do not consider it given the limitations in methodology as discussed in Chapter 6. Also we note that although Pebley, Goldman & Rodríguez (1996)⁵⁴ suggest use of a logistic model, the speed issues discussed in Chapter 4 limit us to a probit link.

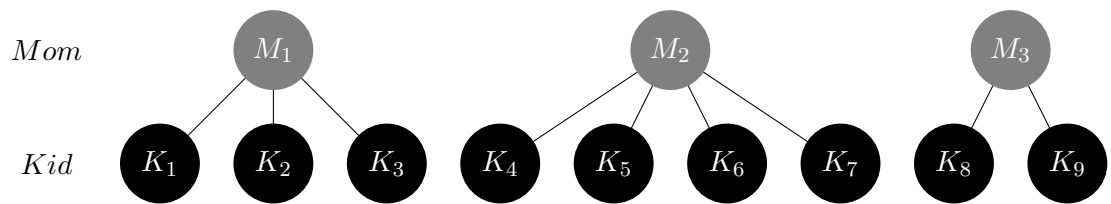


Figure 7.2: A plot showing the two level structure of the data, specifically kids nested under mothers. The first row shows each mother, the second shows each kid of each mother.

The fixed effects included in our model were selected using a generalised additive model selection scheme implemented via the R package “gamsel” (Chouldechova, Hastie & Spinu, 2018^[15]) using an overlap grouped least absolute shrinkage selection operator (Jacob, et al., 2009^[30]), where the most parsimonious model within one standard error from the minimum was chosen (Hastie, Tibshirani & Friedman, 2009^[28]). The package allows including terms as non-linear splines in the model. Error estimates determined using 10 fold cross-validation suggest a good choice for the penalty parameter of a model fitted with least absolute shrinkage selection operator is $\lambda = 2.227$. Although the selected model suggests including primary parental education as a binary variable, given the ordinal nature of education we instead include whether or not the parents finished secondary school for interpretation sake. Linearity between indigenous population percentage in the child’s community and their immunisation status was also confirmed visually using spline fits in the aforementioned package.

We fit the model specified in equation (7.1) with the four following methods and compare the results:

1. **Markov Chain Monte Carlo** via the function `stan()` from the R package “rstan” (Stan Development Team, 2018^[59]), using a Bayesian version of the model specified in equation (7.1) with independent $N(0, 1010)$ distributed diffuse priors for β_0, \dots, β_6 , and a member of the marginally noninformative family of covariance matrix priors described in Huang & Wand (2013)^[29] for the 2×2 covariance matrix in equation (7.2). In the notation of the same paper, the hyperparameters were set to $\nu = 2$ and $A_1 = A_2 = 10^5$. A warm-up of 200000 was used followed by samples of 10000 retained for inference.
2. **Laplace approximation** via the function `glmer()` in the R package “lme4” (Bates, et al., 2018^[5]).

3. **Expectation propagation** as described in Chapter 3 via our function `glmmEP()` in the R package “`glmmEP`” (Wand & Yu, 2020^[69]).
4. **Data cloning** via the function `dclone()` in the R package “`dclone`” (Solymos, 2010^[61]), with 10 clones.

We assume that the Markov Chain Monte Carlo approach gives close to exact results since a quadrature approach is not suitable for more than one random effect. Although other approaches exist, we include a comparison to the data cloning approach since it provides a reliable frequentist inference. The EP point estimates and approximate 95% confidence intervals are given in Table 7.2. With the exception of those involving parental education, each of the parameters is seen to be statistically significantly different from zero. As examples, the EP 95% confidence interval for β_1 of (-1.08, -0.454) indicates a lower prevalence of immunization in communities with higher percentages of indigenous people and the 95% confidence interval for σ_2 of (1.54, 4.35) shows that there is significant heterogeneity in the indigenous percentage effect across the 1595 families.

Figure 7.3 shows comparison of the fixed effects point estimates and approximate 95% confidence intervals of the four methods. Here we use Markov chain Monte Carlo as the gold standard method and apply the assumption that Markov chain Monte Carlo based 95% credible intervals are close to the 95% confidence intervals based on exact maximum likelihood. Although each method provided similar answers regarding significance of parameters, their point estimates and confidence interval coverage were notably different. Only Markov chain Monte Carlo produced a confidence interval for the intercept parameter which was significant. The other methods produced confidence intervals for the intercept parameter with a high proximity to 0, suggesting a lack of evidence to include it. Fixed effect estimates provided by data cloning show considerable bias and reduced standard errors, most notably for indigenous population percentage in the child’s community and age of the child, where point estimates for each parameter were outside the credible interval of the Markov chain Monte Carlo approach. Using data cloning to fit this model took approximately 3hrs on a contemporary laptop, considerably slower than the 12 seconds taken by EP. We also note data cloning has quite a few tuning parameters such as the number of clones and prior hyperparameter values. Given the slow speed of data cloning, it was difficult to assess the sensitivity of parameter choice.

Laplace approximation of the model performed slightly better than the data cloning approach, however was still poor. Given the multilevel model we fit has a low number of

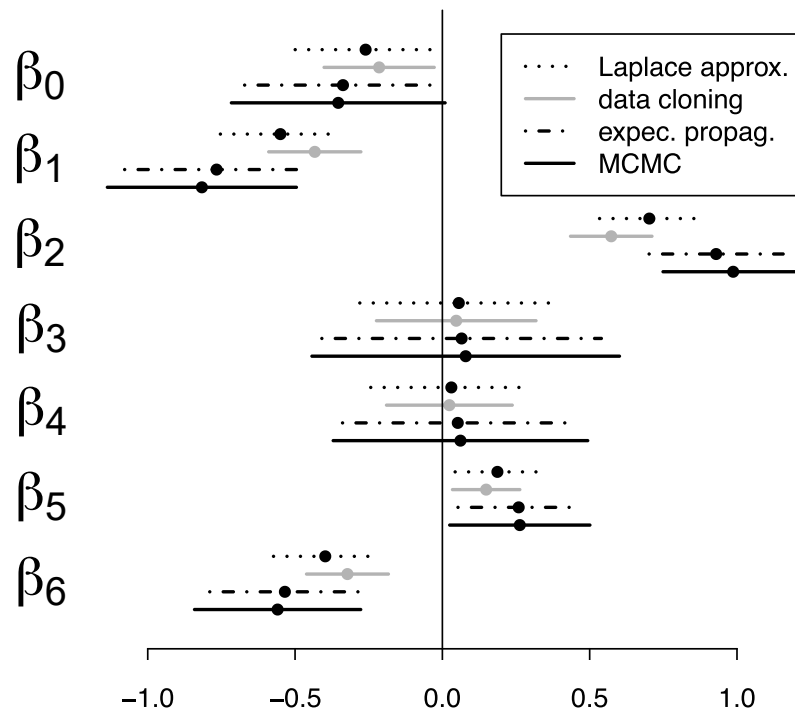


Figure 7.3: Visual comparison of approximate 95% confidence/credible intervals for β_0, \dots, β_6 for three approaches to fitting the probit mixed model equation (7.1) to the Guatemala immunization data. The approaches are Laplace approximation, data cloning, EP and Markov chain Monte Carlo (MCMC) with details given in the text.

observations per group, the breakdown of Laplace approximation's asymptotic properties and a subsequent drop in performance is expected and consistent with literature, as documented in Vonesh (1996).⁶⁵ Additionally, we note that unlike our R function `glmEP()`, the R function `glmer()` does not produce confidence intervals for random effect parameters. As such, although it took 6 seconds (less than half the time of EP) to fit the model with confidence intervals for the fixed effects only, to include confidence intervals for random effects as well took 124 seconds.

Expectation propagation provided consistent point estimates and confidence interval coverage close to Markov chain Monte Carlo estimates whilst being considerably cheaper to compute. For this application, it was clearly the standout method.

Table 7.2: A table of approximate maximum likelihood estimates and corresponding upper and lower 95% confidence intervals given by the EP methodology for the parameters in model equations (7.1) and (7.2).

Parameter	95%C.I. low	Estimate	95%C.I. upp
β_0	-0.6711	-0.3373	-0.0035
β_1	-1.0783	-0.7663	-0.4543
β_2	0.7018	0.9291	1.1565
β_3	-0.4090	0.0653	0.5396
β_4	-0.3388	0.0523	0.4434
β_5	0.0531	0.2591	0.4650
β_6	-0.7895	-0.5345	-0.2795
σ_1	1.1622	1.5370	2.0328
σ_2	1.5407	2.5887	4.3494
ρ	-0.9486	-0.7821	-0.2766

7.3 Modelling donor attendance of the Australian Red Cross Blood Service

Blood donation is currently one of the only ways to collect blood for therapeutic use and as such is vital for healthcare in all nations. Given the often voluntary nature of blood donations, the rate of donor attendance can be quite variable. Such variations have implications for both financial and human resources of the health care system (Boksmati, et al., 2016).^[9] As such, blood collection services seek to understand the factors involved in donor attendance. Although an abundance of work exists in understanding how to maintain steady donations, such as Charbonneau, et al. (2015),^[13] Bagot, et al. (2013)^[2] and Gemelli, et al. (2017),^[20] most have small datasets and fail to provide an analysis using modern statistical tools.

Our dataset provided by the Australian Red Cross Blood Service has over 3 million donation records collected in Queensland from 2015 to 2017. Table 7.3 provides a description of each of the 43 variables recorded. We aim to identify key factors that contribute to blood donors' appointment results (i.e. whether or not the donor will attend their scheduled donation), whilst also including more factors and a larger sample size than other studies. Given the size of the data, computationally efficient methods are required. As such, we compare GLMM fitting using our methodology to Laplace approximation. Furthermore, since the starting values of our EP algorithm depended on estimates from Laplace approximation, fitting this model allows an insight into the

feasibility of the new methodology.

Table 7.3: A table that shows each of the variables in the analysis, explains what they mean, and lists the type of data they are.

Variable name	Description
ID	A multilevel factor variable with each donors' ID number.
attendance	A two-level factor whether the donors either <i>attended</i> or <i>absent</i> their appointment. Didn't attend is the reference.
age	A continuous variable with the donor's age in years.
sex	A two-level factor variable with the donor's gender. Female is the reference.
bloodType	A nine-level factor variable with donor's the blood type. It can be <i>A+</i> , <i>A-</i> , <i>AB-</i> , <i>AB+</i> , <i>B-</i> , <i>B+</i> , <i>O-</i> , <i>O+</i> or <i>Unknown</i> .
pltDon	A count variable with the number of platelet donations the donor has had.
WBDon	A count variable with the number of whole blood donations the donor has had.
plsDon	A count variable with the number of plasma donations the donor has had.
apptTime	A continuous positive variable with time of the appointment, measured from midnight.
apptWkday	A six-level factor variable with weekday of the appointment from 2 to 7, where 2,3,4,5,6,7 are weekdays Monday, Tuesday, Wednesday, Thursday, Friday, Saturday respectively.
apptMth	A 12-level factor variable with the month of the appointment.
apptSeas	A four-level factor variable with season of the appointment.
apptYear	A three-level factor variable with the year of the appointment.
apptPubHolRad	A five-level factor variable with the number of days either before or after the appointment to the nearest public holiday. If a public holiday is not within seven days of an appointment, this is 0. If an appointment is one day from a public holiday, this is 1. If an appointment is two days from a public holiday, this is 2. If an appointment is three day from a public holiday, this is 3, etc.

apptCreDelta	A continuous variable with the difference between date of creating the booking and actual donation date.
apptBookType	A seven-level factor with the type of booking made. <i>In centre</i> bookings are those made after the previous donation. <i>MCC inbound</i> is where the donor calls the centre. <i>MCC inbound</i> is where the centre calls the donor. <i>MCC unknown</i> is a booking made over the phone, but it is unsure who called who. <i>Community relations</i> are bookings made through community outreach. <i>Portal</i> are bookings made using the portal software and <i>Webbooking</i> are those made over the web.
colTypeSched	A four-level factor variable with the type of blood collection procedure scheduled to be conducted. It can be either <i>plasmapheresis</i> , <i>plateletpheresis</i> , <i>whole blood</i> or <i>no blood</i> .
colTypeTaken	A four-level factor variable with the type of blood collection procedure conducted. It can be either <i>plasmapheresis</i> , <i>plateletpheresis</i> , <i>whole blood</i> or <i>no blood</i> taken during appointment.
centreID	A multilevel factor variable with the centre where blood collection was scheduled. There are 25 centres represented by unique codes.
centreType	A two-level factor with the type of centre where the collection of blood was scheduled, i.e. whether it was a <i>static</i> collection room or a <i>mobile</i> blood van.
centreLoc	A two level factor whether the scheduled collection centre was <i>regional</i> or <i>metro</i> . Here <i>regional</i> is the reference category.
lastDefType	A 354 level factor with the donor's last deferral type code.
lastDefLen	A continuous variable with length of time in weeks since the donor last deferred an appointment. This is set to 0 for when the donor has not had a previous deferral.
lastDefEndLen	A continuous variable with the length of time in weeks since the donors last deferral ended. This is set to 0 for when the donor has not had a previous deferral.
inbndSMS	A count variable with the number of SMS to the Red Cross sent by the donor.

<code>inbndPhone</code>	A count variable with the number of phone calls to the Red Cross by the donor.
<code>inbndEmail</code>	A count variable with the number of emails sent to the Red Cross by the donor.
<code>inbndUnknown</code>	A count variable with the number of unknown communications sent to the Red Cross by the donor.
<code>inbndInternet</code>	A count variable with the number of internet communications to the Red Cross by the donor.
<code>outbndPhone</code>	A count variable with the number of phone calls to the donor by the Red Cross.
<code>outbndSMS</code>	A count variable with the number of SMS to the donor sent by the Red Cross.
<code>outbndEmail</code>	A count variable with the number of emails sent to the donor by the Red Cross.
<code>outbndUnknown</code>	A count variable with the number of unknown communications sent to the donor by the Red Cross.
<code>outbndInternet</code>	A count variable with the number of internet communications to the donor by the Red Cross.
<code>lastTTC</code>	A continuous variable with the time it took in the last appointment before the donor was on the couch ready for the phlebotomist to take blood.
<code>lastBleedTime</code>	A continuous variable with the time it took in the last appointment from needle in to needle out. This is linked to collection type.
<code>lastDonTime</code>	A continuous variable the total time it took from when the donor entered the Red Cross to when they finished giving blood, i.e. sum of <code>lastTTC</code> and <code>lastDonTime</code> .
<code>lastDonMultiArm</code>	A two-level factor variable whether or not the donor was punctured in one or both arms. <i>True</i> if one arm and <i>False</i> if both.
<code>lastDAEVVR</code>	A two-level factor variable with the last donor adverse event, specifically whether a vasovagal reaction took place. <i>True</i> if one did happen and <i>False</i> if one did not.

`lastDAEnonVVR` A two-level factor variable with the last donor adverse event, specifically whether a non-vasovagal reaction took place. *True* if one did happen and *False* if one did not.

7.3.1 Data cleaning

The dataset was pooled from many different sources and required cleaning. Donors were free to attend any room as many times as they wished in any location in Queensland, and as such, many donors attended different rooms for their appointments. The underlying structure of the data is quite complicated, involving donations nested in donors, crossed with rooms, where rooms are nested under area (see Figure 7.4). Additionally there is a time effect to the true structure of the data. However, given the limitations of the methodology presented in this thesis, we restrict the data to have one level of nesting; donations nested in donors (see Figure 7.5). We achieve this by filtering the donors so that only visits to their most popular room were retained.

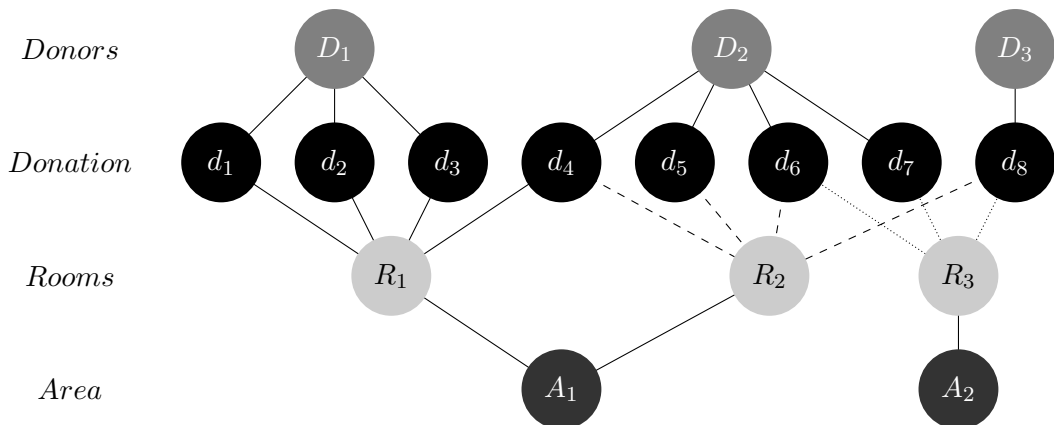


Figure 7.4: A plot showing the two level multiple membership structure of the data, specifically donations nested under donors crossed with rooms, which are nested under area. The first row shows each donor, the second shows each donation of each donor, the third shows which room each donation was done in, and the fourth shows the area of each room. Note the time aspect of the data is not included in this graph.

Some variables were reported as they were used in the collection rooms and as such were not in a format that lent themselves to statistical analysis. `apptTime`, which was in a date time format, was converted to seconds from midnight. The date component was also separated into years.

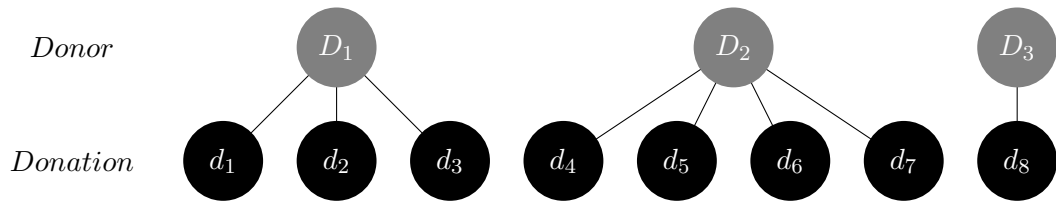


Figure 7.5: A plot showing the one level structure of the data, specifically donations nested under donors. The first row shows each donor, the second shows each donation of each donor.

Patient appointment result was recorded as a six-level factor with outcomes being; attended, blood was taken, cancelled, rescheduled, active or unknown. We converted it to a two-level factor called **attendance**, where donors either attended or were absent, with *absent* as the reference category.

Negative values for `apptCreDelta` were present in the dataset for patients that attended without booking in advance. These negative values were created as the donor's appointment time was recorded after their blood was taken by staff. Since it is not useful for prediction of donors attendance rate, we remove them from the dataset.

There are also a number of missing values present in the data. For the cases of `WBDon`, `pltDon`, `plsDon` and `lastDefLen` missing values may be indicative of first time donors. Given that the analysis provided is primarily a demonstration of methodology, we do not explore this and instead conduct a complete case analysis. A spurious observation where donor's age meant their last deferral length was not possible was also removed from the data.

The counts of contact for each communication method between the Red Cross and the donor were dichotomised from count variables that indicated how many times each form of communication had been conducted on each donor, to a binary variable recording whether there was no communication or at least 1 communication.

Finally there were a number of donors with an extremely high number of observations in the dataset (the highest donor had 2247 bookings). It is speculated that these were automated bookings based on various information flags, that were either not attended or cancelled. Since donors with extremely high numbers of bookings were not uncommon we did not remove them from our dataset.

Table 7.4: A table reporting the shapes and recommended fits of each continuous variable from the R package “gamsel” (Chouldechova, Hastie & Spinu, 2018^[15])

Variable name	Form selected
age	Linear
WBDon	Bend point at 16.
log(plsDon + 1)	Bend points at 1.5, 2.5 and 3.5.
log(pltDon + 1)	Indicator of equalling zero and for positive, bend points at 1.5, 2 and 3.
apptTime	Bend points at 0.45 and 0.53.
apptCreDelta	Bend points at 25 and 50.
log(lastDefLen + 1)	Bend point at 1.2.
log(lastDefEndLen + 1)	Bend points at 1.6, 2.7 and 4.
log(lastTTC + 1)	Indicator of equalling zero and, for positive, bend points at 2, 2.5 and 3.
lastBleedTime	Bend points at 40 and 60.

7.3.2 Modelling continuous variables

We used the R package “gamsel” (Chouldechova, Hastie & Spinu, 2018^[15]) to explore the relationship between each continuous covariate in the dataset and the response variable. Specifically, the R function `gamsel()` was used to automatically select whether variables were linear or non-linear using an overlap grouped least absolute shrinkage selection operator (Tibshirani, 1996^[64]) via the R package “gamsel” (Chouldechova, Hastie & Spinu, 2018^[15]) as in Section 7.2. This showed `plsDon`, `pltDon`, `lastDefLen`, `lastDefEndLen` and `lastTTC` had skewed non-linear relationships with `attendance`. As such we implemented a transform of $\log(x + 1)$ and re-ran the `gamsel()` routine. To account for the non-linearity of some variables we decided to use a broken stick model where the break points were decided by visual inspection of the plots generated by `gamsel()` in conjunction with histograms. The results of this exploration are given in Table 7.4.

Categorical variables were broken into several binary variables each. The reference level of each variable is given in Table 7.3. All continuous variables were also scaled to have mean 0. After cleaning 2430370 observations remained.

7.3.3 Initial model

We first attempted to fit an exploratory model via Laplace approximation. This model included all variables as fixed effects, a random intercept, and random slopes for age. The probit mixed model we fit follows in equation (7.3), where $\mathcal{I}(\mathcal{P}) = 1$ when the logical statement \mathcal{P} is true, $\mathcal{B}(\mathcal{P}; x) = 0$ when the logical statement is false, $\mathcal{B}(\mathcal{P}; x) = x$ when the logical statement is true, and attendance_{ij} is the `attendance` value for the i th donation of the j th donor for $1 \leq i \leq 55657$, $1 \leq j \leq n_i$ and $n_i \in \{1, \dots, 2247\}$. Other variables are defined analogously to `attendanceij`. We assume our bivariate random effects vectors satisfy

$$\begin{bmatrix} u_{0i} \\ u_{1i} \end{bmatrix} \stackrel{\text{ind.}}{\sim} \mathbf{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right). \quad (7.4)$$

Given both the data size and number of fixed effect parameter estimates in the model, the R function `glmer()` in the package “lme4” (Bates, et al., 2018^[5]) was unable to fit the model. As such, we turned to the function `glmmTMB()` (Magnusson, et al., 2019^[39]) from the package “glmmTMB”, which utilises the R package “Template Model Builder” (Kristensen, 2018^[33]) as a backend for model fitting and as such provides a speed advantage for GLMMs with high numbers of fixed and random effects (Bolker, 2019^[7]).

Fitting our model using `glmmTMB()` with all 2430370 observations that remained after cleaning on a high performance computing cluster took approximately 27hrs and used up to 50gb random access memory. We then attempted to fit the same model using our EP approach via a modified version of our R function `glmmEP()` (as in the R package “glmmEP”), with starting values given by the Laplace approximations from `glmmTMB()`. To find the maximum likelihood estimates for the parameters, we first tried using the R function `optim()` to compute a Nelder Mead algorithm with refinements via Broyden Fletcher Goldfarb Shanno algorithm as discussed in Section 3.2. After 2 days of continuous computing time on a high performance computing cluster the Nelder Mead phase of the routine had not completed, and it was decided to try a different method. In particular the unconstrained optimisation by quadratic approximation algorithm was selected for implementation via the function `uobyqa()` in the package “minqa” (Bates, et al., 2015^[6]). The `uobyqa()` approach was significantly faster than the scheme using `optim()` approach, completing in approximately 1 day, and as such we use `uobyqa()` for further investigations.

We tried two methods to calculate the hessian matrix to calculate the confidence intervals for EP. The first method was using Broyden Fletcher Goldfarb Shanno algorithm via the R function `optim()` with `hessian = TRUE`, and the second was using the function `hessian()` from the R package “numDeriv” (Gilbert & Varadhan, 2019^[22]). Both methods failed due to singularities. As such we were not able to obtain confidence intervals for EP. We did not make further attempts to obtain the hessian required since the initial fit was simply exploratory.

Except for the random effects components the difference between point estimates provided by Laplace approximation and EP were minimal. This is most likely due to the high number of donations per donor, which vary between 1 and 2247. The difference in random effects estimates are shown in Table 7.5. Finally we note that majority of the parameter estimates given by Laplace approximation were highly significant, which can be attributed to the large size of the dataset used in the study.

Table 7.5: *Table comparing the random effects point estimates obtained by Laplace approximations and by EP.*

Parameter	Laplace approx.	EP approx.
σ_1	0.4547	0.4544
σ_2	0.2263	0.2261
ρ	0.0538	0.0248

$$\begin{aligned}
 & \mathcal{I}(\text{attendance}_{ij} = \text{attended})|u_{0i}, u_{1i} \stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left\{ \Phi \left\{ \beta_0 + u_{0i} + (\beta_1 + u_{1i}) \text{age}_{ij} \right. \right. \\
 & \quad + \beta_2 \mathcal{I}(\text{sex}_{ij} = \text{male}) + \beta_3 \mathcal{I}(\text{bloodType}_{ij} = A+) + \beta_4 \mathcal{I}(\text{bloodType}_{ij} = A-) \\
 & + \beta_5 \mathcal{I}(\text{bloodType}_{ij} = AB-) + \beta_6 \mathcal{I}(\text{bloodType}_{ij} = AB+) + \beta_7 \mathcal{I}(\text{bloodType}_{ij} = B+) \\
 & + \beta_8 \mathcal{I}(\text{bloodType}_{ij} = B-) + \beta_9 \mathcal{I}(\text{bloodType}_{ij} = O+) + \beta_{10} \mathcal{I}(\text{bloodType}_{ij} = O-) \\
 & \quad + \beta_{11} \text{WBDon}_{ij} + \beta_{12} \mathcal{B}(\text{WBDon}_{ij} > 16) + \beta_{13} \log(\text{plsDon} + 1) \\
 & \quad + \beta_{14} \mathcal{B}(\log(\text{plsDon} + 1) > 1.5) + \beta_{15} \mathcal{B}(\log(\text{plsDon} + 1) > 2.5) \\
 & \quad + \beta_{16} \mathcal{B}(\log(\text{plsDon} + 1) > 3.5) + \beta_{17} \log(\text{pltDon} + 1) \\
 & \quad + \beta_{18} \mathcal{I}(\log(\text{pltDon} + 1) = 0) + \beta_{19} \mathcal{B}(\log(\text{pltDon} + 1) > 1.5) \\
 & + \beta_{20} \mathcal{B}(\log(\text{pltDon} + 1) > 2) + \beta_{21} \mathcal{B}(\log(\text{pltDon} + 1) > 3) + \beta_{22} \text{apptTime} \\
 & + \beta_{23} \mathcal{B}(\text{apptTime} > 0.45) + \beta_{24} \mathcal{B}(\text{apptTime} > 0.53) + \beta_{25} \mathcal{I}(\text{apptWkday} = 3) \\
 & + \beta_{26} \mathcal{I}(\text{apptWkday} = 4) + \beta_{27} \mathcal{I}(\text{apptWkday} = 5) + \beta_{28} \mathcal{I}(\text{apptWkday} = 6) \\
 & + \beta_{29} \mathcal{I}(\text{apptWkday} = 7) + \beta_{30} \mathcal{I}(\text{apptMth} = 2) + \beta_{31} \mathcal{I}(\text{apptMth} = 3) \\
 & + \beta_{32} \mathcal{I}(\text{apptMth} = 4) + \beta_{33} \mathcal{I}(\text{apptMth} = 5) + \beta_{34} \mathcal{I}(\text{apptMth} = 6) \\
 & + \beta_{35} \mathcal{I}(\text{apptMth} = 7) + \beta_{36} \mathcal{I}(\text{apptMth} = 8) + \beta_{37} \mathcal{I}(\text{apptMth} = 9) \\
 & + \beta_{38} \mathcal{I}(\text{apptMth} = 10) + \beta_{39} \mathcal{I}(\text{apptMth} = 11) + \beta_{40} \mathcal{I}(\text{apptMth} = 12) \\
 & + \beta_{41} \mathcal{I}(\text{apptYear} = 2016) + \beta_{42} \mathcal{I}(\text{apptYear} = 2017) + \beta_{43} \mathcal{I}(\text{apptPubHolRad} = 1) \\
 & + \beta_{44} \mathcal{I}(\text{apptPubHolRad} = 2) + \beta_{45} \mathcal{I}(\text{apptPubHolRad} = 3) + \beta_{46} \mathcal{I}(\text{apptPubHolRad} = 4) \\
 & + \beta_{47} \text{apptCreDelta} + \beta_{48} \mathcal{B}(\text{apptCreDelta} > 25) + \beta_{49} \mathcal{B}(\text{apptCreDelta} > 50) \\
 & + \beta_{50} \mathcal{I}(\text{apptBookType} = \text{community relations}) + \beta_{51} \mathcal{I}(\text{apptBookType} = \text{in centre}) \\
 & + \beta_{52} \mathcal{I}(\text{apptBookType} = \text{NCC inbound}) + \beta_{53} \mathcal{I}(\text{apptBookType} = \text{NCC outbound}) \\
 & + \beta_{54} \mathcal{I}(\text{apptBookType} = \text{portal}) + \beta_{55} \mathcal{I}(\text{apptBookType} = \text{web booking}) \\
 & + \beta_{56} \mathcal{I}(\text{colTypeSched} = \text{no blood}) + \beta_{57} \mathcal{I}(\text{colTypeSched} = \text{plasmapheresis}) \\
 & + \beta_{58} \mathcal{I}(\text{colTypeSched} = \text{plateletpheresis}) + \beta_{59} \mathcal{I}(\text{colTypeTaken} = \text{no blood}) \\
 & + \beta_{60} \mathcal{I}(\text{colTypeTaken} = \text{plasmapheresis}) + \beta_{61} \mathcal{I}(\text{colTypeTaken} = \text{plateletpheresis}) \\
 & + \beta_{62} \mathcal{I}(\text{centreType} = \text{static}) + \beta_{63} \mathcal{I}(\text{centreType} = \text{metro}) + \beta_{64} \log(\text{lastDefLen} + 1) \\
 & + \beta_{65} \mathcal{B}(\log(\text{lastDefLen} + 1) > 1.2) + \beta_{66} \log(\text{lastDefEndLen} + 1) \\
 & + \beta_{67} \mathcal{B}(\log(\text{lastDefEndLen} + 1) > 1.6) + \beta_{68} \mathcal{B}(\log(\text{lastDefEndLen} + 1) > 2.7) \\
 & + \beta_{69} \mathcal{B}(\log(\text{lastDefEndLen} + 1) > 4) + \beta_{70} \mathcal{I}(\text{inbndPhone} \leq 1) + \beta_{71} \mathcal{I}(\text{inbndSMS} \leq 1) \\
 & + \beta_{72} \mathcal{I}(\text{inbndEmail} \leq 1) + \beta_{73} \mathcal{I}(\text{inbndInternet} \leq 1) + \beta_{74} \mathcal{I}(\text{outbndPhone} \leq 1) \\
 & + \beta_{75} \mathcal{I}(\text{outbndSMS} \leq 1) + \beta_{76} \mathcal{I}(\text{outbndEmail} \leq 1) + \beta_{77} \mathcal{I}(\text{outbndInternet} \leq 1) \\
 & + \beta_{78} \mathcal{I}(\text{outbndLetter} \leq 1) + \beta_{79} \log(\text{lastTTC} + 1) + \beta_{80} \mathcal{I}(\log(\text{lastTTC} + 1) = 0) \\
 & + \beta_{81} \mathcal{B}(\log(\text{lastTTC} + 1) > 2) + \beta_{82} \mathcal{B}(\log(\text{lastTTC} + 1) > 2.5) \\
 & + \beta_{83} \mathcal{B}(\log(\text{lastTTC} + 1) > 3) + \beta_{84} \text{lastBleedTime} + \beta_{85} \mathcal{B}(\text{lastBleedTime} > 40) \\
 & + \beta_{86} \mathcal{B}(\text{lastBleedTime} > 60) + \beta_{87} \mathcal{I}(\text{lastDonMultiArm} = \text{True}) \\
 & \left. + \beta_{88} \mathcal{I}(\text{lastDAEVVR}_{ij} = \text{True}) + \beta_{89} \mathcal{I}(\text{lastDAEnonVVR}_{ij} = \text{True}) \right\}, \quad (7.3)
 \end{aligned}$$

7.3.4 Second model

With insights from the previous model, we aimed to improve model fit and obtain confidence intervals for the EP approach. We conducted variable reduction of the fixed effects using least absolute shrinkage selection operator (Tibshirani, 1996^[64]) via the R package “glmnet” (Friedman, J., et al., 2020^[17]). A plot of the error estimates given the penalty parameter λ determined via 10 fold cross-validation is shown in Figure [7.8](#). Although grounds for using the “one-standard-error” rule are established (Hastie, Tibshirani & Friedman, 2009^[28]) we instead used $\lambda = \exp(-6)$ as suggested by Figure [7.8](#), since it provides a simpler model without an excessive increase in the standard error. Again we included a random intercept and slope for age. The probit mixed model we fit follows in equation [\(7.5\)](#), where attendance_{ij} , $\mathcal{B}(\mathcal{P}; x)$ and $\mathcal{I}(\mathcal{P})$ are as previously

$$\begin{aligned}
 \mathcal{I}(\text{attendance}_{ij} = \text{attended})|u_{0i}, u_{1i} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left\{ \Phi \left\{ \beta_0 + u_{0i} + (\beta_1 + u_{1i}) \text{age}_{ij} \right. \right. \\
 &+ \beta_2 \mathcal{I}(\text{sex}_{ij} = \text{male}) + \beta_3 \mathcal{I}(\text{bloodType}_{ij} = A+) + \beta_4 \mathcal{I}(\text{bloodType}_{ij} = B+) \\
 &\quad + \beta_5 \text{WBDon}_{ij} + \beta_6 \mathcal{B}(\text{WBDon}_{ij} > 16) + \beta_7 \log(\text{plsDon} + 1) \\
 &\quad + \beta_8 \mathcal{B}(\log(\text{plsDon} + 1) > 1.5 \log) + \beta_9 \mathcal{B}(\log(\text{plsDon} + 1) > 2.5) \\
 &\quad + \beta_{10} \mathcal{B}(\log(\text{plsDon} + 1) > 3.5) + \beta_{11} \mathcal{I}(\log(\text{pltdon}_{ij} + 1) = 0) \\
 &+ \beta_{12} \text{apptTime} + \beta_{13} \mathcal{B}(\text{apptTime} > 0.45) + \beta_{14} \mathcal{B}(\text{apptTime} > 0.53) \\
 &\quad + \beta_{15} \mathcal{I}(\text{apptWkday}_{ij} = 3) + \beta_{16} \mathcal{I}(\text{apptWkday}_{ij} = 7) \\
 &\quad + \beta_{17} \mathcal{I}(\text{apptMth}_{ij} = 11) + \beta_{18} \mathcal{I}(\text{apptMth}_{ij} = 12) \\
 &\quad + \beta_{19} \mathcal{I}(\text{apptYear}_{ij} = 2017) + \beta_{20} \text{apptCreDelta} \\
 &\quad + \beta_{21} \mathcal{B}(\text{apptCreDelta} > 25) + \beta_{22} \mathcal{B}(\text{apptCreDelta} > 50) \\
 &\quad + \beta_{23} \mathcal{I}(\text{apptBookType}_{ij} = \text{MCC unknown}) \\
 &\quad + \beta_{24} \mathcal{I}(\text{colTypeSched}_{ij} = \text{plasmapheresis collection}) \\
 &\quad + \beta_{25} \mathcal{I}(\text{colTypeTaken}_{ij} = \text{plasmapheresis collection}) \\
 &+ \beta_{26} \log(\text{lastDefLen} + 1) + \beta_{27} \mathcal{B}(\log(\text{lastDefLen} + 1) > 1.2) \\
 &\quad + \beta_{28} \log(\text{lastDefEndLen} + 1) + \beta_{29} \mathcal{I}(\text{inbndPhone} > 1) \\
 &\quad + \beta_{30} \mathcal{I}(\text{outbndPhone} > 1) + \beta_{31} \mathcal{I}(\text{outbndSMS} > 1) \\
 &\quad \left. \left. + \beta_{32} \mathcal{I}(\text{outbndEmail} \geq 1) + \beta_{33} \mathcal{I}(\text{outbndLetter} \geq 1) \right\} \right\}, \tag{7.5}
 \end{aligned}$$

defined. Other variables are defined analogously to attendance_{ij} . We assumed our bivariate random effects vectors satisfy equation [\(7.4\)](#).

As before we first obtained fits using Laplace approximation via `glmmTMB()` and used

them as starting values for the EP algorithm. Again both models were extremely heavy computationally, with each method taking approximately one day to fit. A difference is noticeable between the estimates of β_{20} , where Laplace approximation provided much tighter confidence intervals than EP. This is a trend that is consistent across almost all parameters estimated, however much like the first fit, the fixed effect confidence intervals and point estimates for Laplace approximations are extremely close to those provided by EP. Figure 7.7 shows a comparison between Laplace approximation and EP. The random effects estimates were also extremely close between the two methods, as shown in Table 7.6. Differences between the models random effects predictions are somewhat visible in Figure 7.6, which shows best predictions of the random intercepts plotted against best predictions of slopes for 24 randomly chosen donors. While these differences are quite small, they are certainly notable for certain donors, such as donor 10, 18 and 29. As a whole however, there is not much separating the fit of each model. We present the point estimates with upper and lower confidence intervals for each parameter in Table 7.7.

7.3.4.1 Results of model fit

All of the parameters estimated in the model were statistically different from zero, except the parameter β_{10} (95% confidence intervals of $(-0.004, 0.0448)$) which was included to account for the nonlinear relationship of the number of plasma donations and attendance. Due to the large dataset, all the fixed effects parameters estimated have extremely tight confidence intervals, with the distance between upper and lower intervals ranging from 0.0035 to 0.0780. However, many variables have coefficients that are small in magnitude and as such although they are significant statistically their impact is minimal in reality.

Age (β_1 with 95% confidence interval $(0.0139, 0.0148)$) increased the probability of donation attendance, as did being male (β_2 with 95% confidence interval $(0.0626, 0.074)$) over being female. Blood types $A+$ and $B+$ (β_3 with 95% confidence interval $(0.0068, 0.0187)$ and β_4 with 95% confidence interval $(0.0159, 0.0293)$ respectively) increased the probability of donation attendance compared to all other blood types, although $B+$ had a larger effect on donation attendance than $A+$.

The number of whole blood donations per donor had a varying effect on donor attendance, in both magnitude and direction. For low values it decreased donation attendance (β_5 with 95% confidence interval $(-0.0046, -0.0023)$), however began increasing it after 16 donations (β_6 with 95% confidence interval $(0.0128, 0.0177)$). The relationship between the number of plasma donations per donor and attendance also varies in

magnitude, but is always negative (β_7 with 95% confidence interval $(-0.0849, -0.0676)$). It weakens between three (β_8 with 95% confidence interval $(0.0561, 0.0919)$) and 11 plasma donations (β_9 with 95% confidence interval $(-0.0575, -0.0158)$), before returning beyond original levels after 32 plasma donations (β_{10} with 95% confidence interval $(-0.004, 0.0448)$). Donors with no platelet donations (β_{11} with 95% confidence interval $(0.0017, 0.013)$) had increased attendance compared to those with platelet donations.

Donor last deferral length has a non-linear relationship with attendance, initially decreasing attendance up to approximately 2.5 days (β_{26} with 95% confidence interval $(-0.0243, -0.012)$), then increasing after this point (β_{27} with 95% confidence interval $(0.0684, 0.0869)$). Time since the last deferral ended was positively linked with donation attendance probability (β_{28} with 95% confidence interval $(0.0005, 0.0041)$).

Both scheduled and collected plasmapheresis donations (β_{24} with 95% confidence interval $(0.0276, 0.0331)$ and β_{25} with 95% confidence interval $(0.0052, 0.0103)$ respectively) were linked with increasing donation attendance.

Appointment time (β_{12} with 95% confidence interval $(-0.8924, -0.7444)$) was linked with decreasing attendance. Much like previous parameters, it has a non-linear relationship with attendance. After appointment time becomes greater than 0.45 (i.e. after approximately 10:45am) (β_{13} with 95% confidence interval $(0.3564, 0.6778)$) the negative effect on probability of attendance decreases, with a further decrease occurring for appointment times greater than 0.53 (i.e. after approximately 12:45am) (β_{14} with 95% confidence interval $(0.5093, 0.7761)$). Note there is an even spread of attendance times between .3 and 0.8 (i.e., 7:12am to 7:12pm), even though it can range from 0 to 1 for all patients, where 0 is 0hrs from midnight and 1 is 24hrs from midnight (i.e., each hour of the day adds $1/24$).

Both the appointment month and day affect attendance probability. Appointments on Wednesdays and Saturdays (β_{15} with 95% confidence interval $(0.0139, 0.0184)$ and β_{16} with 95% confidence interval $(0.0059, 0.0116)$ respectively) were positively linked with donation attendance compared to other days in the week, although the effect size of appointments on Saturdays is half that on Wednesdays. The same trend is true for November and December appointments (β_{17} with 95% confidence interval $(0.008, 0.0115)$ and β_{18} with 95% confidence interval $(0.019, 0.0225)$ respectively), where both months were positively linked with donation attendance compared to other months. We also note donors in 2017 had an increased probability of attending appointments (β_{19} with 95% confidence interval $(0.1095, 0.1142)$).

The length of time between booking and attending the appointment (β_{20} with 95% confidence interval $(-0.0298, -0.0292)$) was negatively linked with donation attendance. As the length of time between booking and attending the appointment becomes greater than 25 days length (β_{21} with 95% confidence interval $(0.0266, 0.0276)$) negative effect becomes minimised, before increasing again as the length of time between booking and attending the appointment becomes greater than 50 days length (β_{22} with 95% confidence interval $(-0.0019, -0.001)$).

Appointments with booking type MCC unknown (β_{23} with 95% confidence interval $(0.2577, 0.2619)$) were the most positively linked variable with probability of donation attendance.

We found more than one text message from the Red Cross to the donor (β_{31} with 95% confidence interval $(0.0846, 0.0887)$) led to the largest increase in probability of donation attendance of any communication method. Letters and emails sent from the Red Cross to the donor lead to smaller increases in attendance probability (β_{33} with 95% confidence interval $(0.0123, 0.0162)$ and β_{32} with 95% confidence interval $(0.0065, 0.0106)$ respectively). Phone calls from the Red Cross to the donor and from the donor to the Red Cross were both linked to a decrease in the probability of appointment attendance (β_{30} with 95% confidence interval $(-0.0355, -0.0313)$ and β_{29} with 95% confidence interval $(-0.0137, -0.0096)$ respectively).

The random intercept parameter (σ_1 with 95% confidence interval $(0.4826, 0.4952)$) shows attendance probability of donors varies, while the random slope (σ_2 with 95% confidence interval $(0.2352, 0.2573)$) suggests the effect of age varies even less between patients. The correlations of the random intercept and slope (ρ with 95% confidence interval $(0.0017, 0.0591)$) are minimal and are likely to be negligible in reality.

In summary, from our model we are able to identify the following useful trends:

- The effect of age is dependent on the donor although in general older patients are more likely to attend their appointments, as are males compared to females. These results are supported by Gemelli, et al. (2017)^[20] and Charbonneau, et al. (2015)^[13] who shows older patients and males are more frequent donors. The variability of the effect of age may be caused by older patients being more experienced donors, in addition to other behavioural effects that occur with aging. Charbonneau, et al. (2015)^[13] suggest females may stop attending due pregnancy, or may stop miss donations based on their menstrual cycle.
- Donors with blood type $B+$ are associated with attendance more than other blood

groups. This contradicts Gemelli, et al. (2017)^[20] who found negative blood groups more common than positive groups.

- High numbers of previous plasma donations are linked to decreasing attendance, while the number of platelet and whole blood donations have negligible effect. Charbonneau, et al. (2015)^[13] also supports the idea that the number of plasma donations has an affect on donation attendance. This may be because there is a large gap required between whole blood donations (nearly two months), where as plasma donations can be once a month. Although platelet donations can be once every seven days, donations can take up to 3 hours which may deter patients from returning.
- A long last deferral length is favourable for attendance, while time since the last deferral does not have any effect. Although the reasoning behind this is not entirely clear the result aligns with the findings in Spekman, et. al. (2019)^[40]
- Donors with scheduled plasmapheresis donations are most likely to attend compared to other appointments. This might be because the minimum time between appointments is the shortest of all the appointment types.
- Early and late appointment times seem to have opposite effects on attendance, where early appointments decrease attendance and late appointments increase attendance. This may be due to patients work commitments, or having greater confidence to donate in later parts of the day after eating several meals.
- Wednesday and December are associated more with attendance than other weekdays and months respectively. Although it is difficult to intuit why Wednesday is associated in increasing donor attendance, it does make sense that people may have spare time over the holiday period during December, or may be influenced by the festive season.
- The booking type MCC unknown seems to give the largest improvement in donor attendance, however there is not much intuition behind this.
- More than one outbound text message led to the largest increase in attendance of the communication methods, followed by outbound letters. Gemelli et. al. (2018)^[21] supports that SMS increased the odds of donors attendance. Outbound and inbound phone calls decreased attendance, possibly due to the large number of cancellation calls, deferrals or missing patient calls.

Table 7.6: Table comparing the random effects point estimates obtained by Laplace approximations and by EP for the model specified in Equation (7.5).

Parameter	Laplace approx.	EP approx.
σ_1	0.4876	0.4889
σ_2	0.2441	0.2460
ρ	0.0300	0.0304

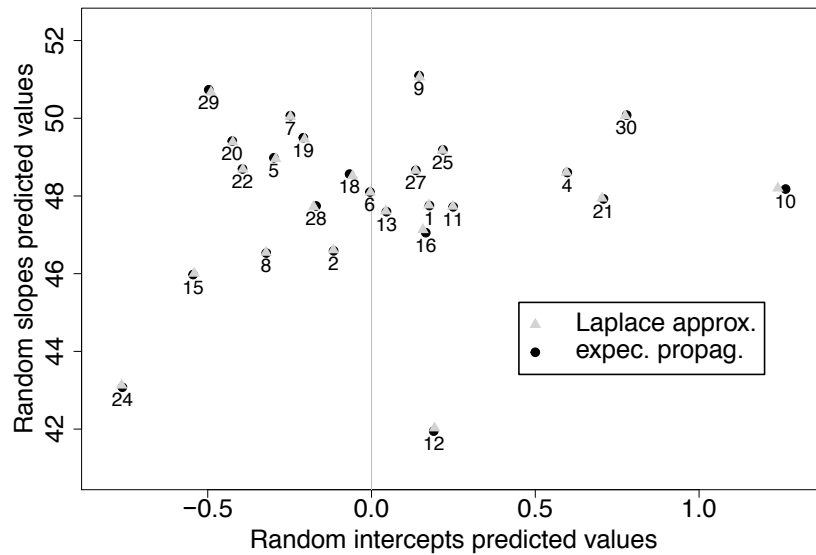


Figure 7.6: Comparison plot for best predictions of the random intercepts and slopes for 24 donors between Laplace approximation and EP.

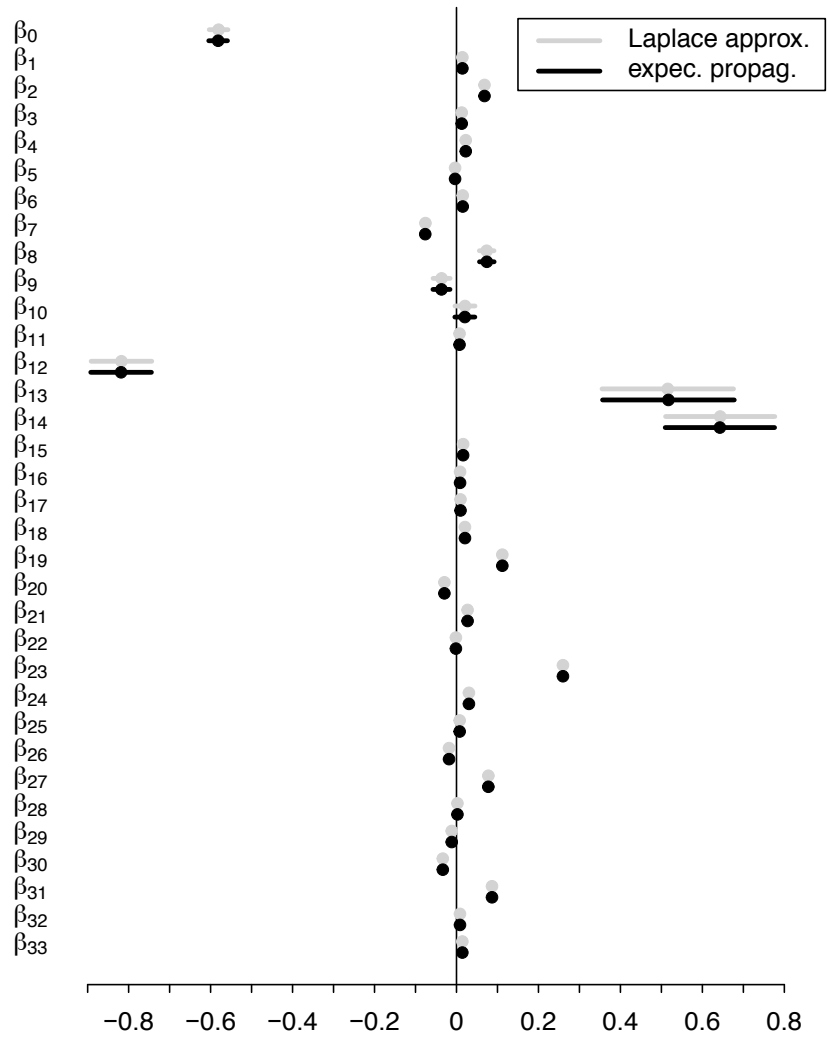


Figure 7.7: Comparison plot of the confidence intervals and point estimates obtained with Laplace approximations and EP for the fixed effect parameters corresponding to the model specified in equation (7.5).

7.4 Appendix

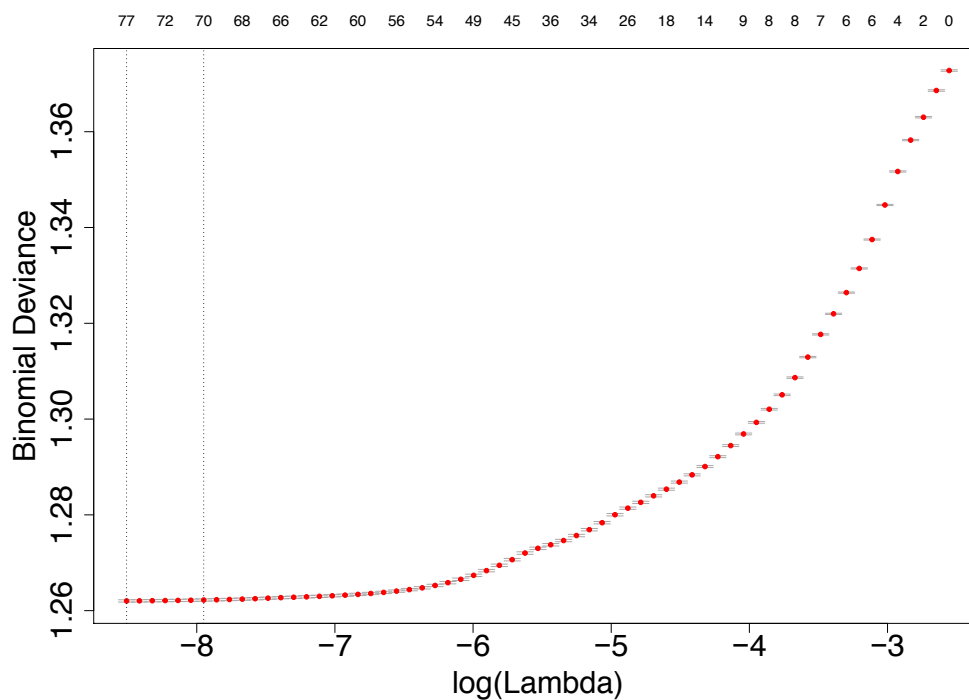


Figure 7.8: Plot of error estimates given the penalty parameter λ determined via 10 fold cross-validation. The values shown horizontally on top of the plot are the number of variables included in the model.

Table 7.7: Table listing the fixed effects point estimates obtained by EP according to Equation (7.5). Note that they are the same as those by EP.

Parameter	95% C.I. low	Estimate	95% C.I. high
β_0	-0.6045	-0.5816	-0.5588
β_1	0.0139	0.0144	0.0148
β_2	0.0626	0.0683	0.074
β_3	0.0068	0.0128	0.0187
β_4	0.0159	0.0226	0.0293
β_5	-0.0046	-0.0034	-0.0023
β_6	0.0128	0.0152	0.0177
β_7	-0.0849	-0.0762	-0.0676
β_8	0.0561	0.074	0.0919
β_9	-0.0575	-0.0366	-0.0158
β_{10}	-0.004	0.0204	0.0448
β_{11}	0.0017	0.0074	0.013

Parameter	95% C.I. low	Estimate	95% C.I. high
β_{12}	-0.8924	-0.8184	-0.7444
β_{13}	0.3564	0.5171	0.6778
β_{14}	0.5093	0.6427	0.7761
β_{15}	0.0139	0.0162	0.0184
β_{16}	0.0059	0.0087	0.0116
β_{17}	0.008	0.0098	0.0115
β_{18}	0.019	0.0207	0.0225
β_{19}	0.1095	0.1119	0.1142
β_{20}	-0.0298	-0.0295	-0.0292
β_{21}	0.0266	0.0271	0.0276
β_{22}	-0.0019	-0.0014	-0.001
β_{23}	0.2577	0.2598	0.2619
β_{24}	0.0276	0.0303	0.0331
β_{25}	0.0052	0.0077	0.0103
β_{26}	-0.0243	-0.0182	-0.012
β_{27}	0.0684	0.0777	0.0869
β_{28}	5e-04	0.0023	0.0041
β_{29}	-0.0137	-0.0116	-0.0096
β_{30}	-0.0355	-0.0334	-0.0313
β_{31}	0.0846	0.0867	0.0887
β_{32}	0.0065	0.0085	0.0106
β_{33}	0.0123	0.0143	0.0162
σ_1	0.4826	0.4889	0.4952
σ_2	0.2352	0.246	0.2573
ρ_{12}	0.0017	0.0304	0.0591

Chapter 8

Discussion and conclusion

This thesis aimed to add to the body of methodology used for frequentist inference of GLMMs by using the typically Bayesian idea of EP. We begun developing our methodology in Chapter 2 on a random intercept only probit model, as well as exploring possible solutions for computing confidence intervals. We explained how EP is used in our model, how we obtained starting values for the EP algorithm, and that the EP approximation of the likelihood surface is very similar around the maximum to the exact likelihood. Three possible solutions for obtaining confidence intervals were explored and we considered how they performed via a simulation study. For this study, approach II of EP suffered significantly in terms of computational performance. Keeping in mind the eventual extension to higher dimensional random effects, it was most practical to use a quasi-Newton optimisation approach due to the algebraic cost of implementing the other approaches. However, assuming the random effects dimension remains constant, approach I is note worthy for its potential computational speed given it requires only a bisection search. Best prediction theory for the univariate model was also covered.

The generalised case of the probit model from Chapter 2 is explored in Chapter 3. Specifically, in this chapter we developed methodology to handle any number of fixed and random effects for one level models. We first explained how to obtain the EP likelihood approximation and the starting values required for the algorithm, before covering how to compute point estimates and confidence intervals using the quasi-Newton method selected in Chapter 2, as well as best predictors of the random effects. Our simulation study comparing the main approaches (Laplace approximation and Gauss-Hermite quadrature) and as well as a Markov chain Monte Carlo based method proved our EP approach is more consistent and accurate than other approaches in a variety of settings. In practice our methodology is hindered computationally given it is dependent

on Laplace approximations for starting values. Additional computational performance may be found by implementing more advanced optimisation algorithms over a limited search area. Although Greene (2012)^[23] suggests against the use of BFGS optimisation, simulation studies throughout the thesis provide evidence that it is suitable, with good empirical coverage provided for maximum likelihood estimates.

We next explored methods for handling logistic models with one level of nesting in Chapter 4. The projections for the logistic model were difficult to solve since closed form solutions do not exist for them. In the univariate case we considered the use of a piecewise approximation, which compared well to the exact likelihood. However, the piecewise approximation did not extend well to the multivariate setting, and as such we explored solving the projections using quadrature. We proved a series of results which allowed the expression of the required multivariate integrals as univariate integrals. Additionally we showed how to express these integrals in a stable manner for computation. The quadrature method compared well to Laplace approximation with regards to accuracy, however had poor computational speed. This speed problem was caused by the use of trapezoidal quadrature. Although it was possible to solve these integrals using Gauss-Hermite quadrature, it was not possible to gauge the level of accuracy the quadrature provided. To our knowledge there is not a quadrature rule that allows approximation of the integral occurring in this chapter within an error bound. We left this as an open problem for future research.

Following the work of Chapter 4, in Chapter 5 we explored Poisson and negative-binomial models for count data. As with the logistic model, the projections required for both count models were not available with closed form solutions. As such we utilised similar results to express the integrals required in a stable and computationally efficient form. In the Poisson model we showed that for the settings of the simulation study, EP provided marginally superior coverage to Laplace approximations, however Gauss-Hermite quadrature appeared the closest to 95% empirical coverage. Similar findings were presented in the negative binomial model, where our EP method provided marginally better results than Laplace approximations and Gauss-Hermite quadrature. We provided confidence intervals for the shape parameter in the negative binomial model, which to our knowledge is not available elsewhere in literature. Given the time limitations of the thesis and complexity of the coding involved we did not present results of multivariate studies for count data. Additionally, we note that although we are able to give confidence intervals for the shape parameter for the negative binomial model, further research into the implications of doing so are required. Finally a careful analysis of the negative binomial model implementation provided in this thesis may be beneficial to providing better numerical stability as the shape parameter approaches infinity.

We investigated the extensions to higher level models in Chapter 6. We explained how crossed random effects and two level models can be framed and implemented using the message passing framework shown in the previous chapters, with only minor revisions to the algebra. Additionally we showed the simple alterations for calculation of confidence intervals. Simulation results showed the performance of the EP approach was similar to that of Laplace approximations, although given the challenging dataset no method performed particularly well. Due to time limitations, we were not able to conduct a more indepth exploration of EP in these models, and leave this as an open problem. Although we did not explicitly show it, the extensions shown in Chapters 3-4 should help give an idea of the steps required to accomodate for crossed random effects models with different link functions and responses types.

In Chapter 7 we applied methodology developed in Chapter 3 to two datasets. In the first dataset, we showed how EP is useful for data that is difficult to model with current methods. Specifically, we showed that for data with a low number of observations per grouping variable, EP is not only as accurate as MCMC, but also provides a good compromise in terms of speed. In the second dataset, where the number of observations per grouping variable was high, we showed that there is only a very minor difference between EP and Laplace approximation. Given the computational cost of implementing EP for this dataset, Laplace approximation was favourable. In this second analysis we also provided useful information on the practicalities of fitting large GLMMs particularly from the software standpoint, where not only do we suggest useful alternatives in the case where the ubiquitous “lme4” package breaks, but also contribute our own software package “glmmEP” (Wand & Yu, 2020⁶⁹). Although we provide some interesting and novel findings from our analysis, a further exploration of the dataset considering things such as missing values and the spacio-temporal structure of the data could be useful for further inference. Additionally, it is possible to obtain starting values for EP by using a sample of the whole dataset, which may lead to computational savings. Further research into this may be of great use for applications of the EP methodology to large datasets. Finally, the large number of events for some patients is relevant to the underlying nature of the dataset and should not be omitted. In no major way would Figure 7.6 be effected by these observations other than both EP and Laplace approximation having similar estimates of fixed and random effects for those patients. The author notes that covergence issues experienced in the first model fit for the Australian Red Cross blood data was caused by highly correlated fixed effects.

In conclusion, this thesis has led to the development of novel methodology for one level binary and count GLMMs, as well as binary crossed random effects GLMMs. We created an efficient message passing framework for researchers to expand and develop

upon, allowing EP, an idea typically implemented in Bayesian settings, to be used for frequentist inference. Additionally we contributed the R package “glmmEP” (Wand & Yu, 2020^[69]) and showed our method provides consistent and accurate parameter estimates for GLMMs, often with a clear improvements in empirical coverage variance parameters. Our work contributes to the underdeveloped frequentist setting for datasets with low numbers of observations per group.

References

- [1] Azzalini, A. (2020). `sn`: The Skew-Normal and Related Distributions Such as the Skew-t (Version 1.5-5). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/sn/index.html>
- [2] Bagot, K. L., Bove, L. L., Masser, B. M., Bednall, T. C., & Buzza, M. (2013). Perceived deterrents to being a plasmapheresis donor in a voluntary, nonremunerated environment. *Transfusion*, 53(5), 1108–1119.
- [3] Barthelme, S. (2016). Part 1: The Expectation-Propagation algorithm, lecture recording, Week 5: “Bayesian statistics and algorithms”, Centre International de Rencontres Mathématiques, delivered 2 March 2016.
- [4] Bates, D., Maechler, M., & Bolker, B. (2019). `mlmRev`: Examples from Multilevel Modelling Software Review (Version 1.0-7). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/mlmRev/index.html>
- [5] Bates, D., Maechler, M., Bolker, B., Walker, S., & Christensen, R. H. B. (2018). `lme4`: Linear Mixed-Effects Models using ‘Eigen’ and S4 (Version 1.1-17). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/lme4/lme4.pdf>
- [6] Bates, D., Mullen, K. M., Nash, J. C., & Varadhan, R. (2014). `minqa`: Derivative-free optimization algorithms by quadratic approximation (Version 1.2.4). Retrieved from <https://cran.r-project.org/web/packages/minqa/index.html>
- [7] Bolker, B. (2019). Getting started with the `glmmTMB` package. <https://cran.r-project.org/web/packages/glmmTMB/vignettes/glmmTMB.pdf>
- [8] Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H., & White, J. S. S. (2009). Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in Ecology & Evolution*, 24(3), 127–135.
- [9] Boksmati, N., Butler-Henderson, K., Anderson, K., & Sahama, T. (2016). The Effectiveness of SMS Reminders on Appointment Attendance: a Meta-Analysis. *Journal of Medical Systems*, 40(4), 90
- [10] Borenstein, M., Hedges, L. V., Higgins, J. P. T., & Rothstein, H. R. (2009). *Introduction to Meta-Analysis*. West Sussex: John Wiley & Sons Ltd.
- [11] Breslow, N. E., & Clayton, D. G. (1993). Approximate Inference in Generalized Linear Mixed Models. *Journal of the American Statistical Association*, 88(421), 9–25.

- [12] Carlin, B. P., & Gelfand, A. E. (1991). An for iterative Monte Carlo method nonconjugate Bayesian analysis. *Statistics and Computing*, 1, 119–128.
- [13] Charbonneau, J., Cloutier, M. S. & Carrier, É. (2015). Why Do Blood Donors Lapse or Reduce Their Donation’s Frequency? *Transfusion Medicine Reviews*, 30(1), 1–5.
- [14] Chesher, A. (1997). Non-normal variation and regression to the mean. *Stat Methods Med Res.* 6(2), 147–66.
- [15] Chouldechova, A., Hastie, T., & Spinu, V. (2018). *gamsel: Fit Regularization Path for Generalized Additive Models*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/gamsel/gamsel.pdf>
- [16] Collins, D. (2008). The performance of estimation methods for generalized linear mixed models. (Doctor of Philosophy), University of Wollongong, Wollongong.
- [17] Friedman, J., Hastie, T., Tibshirani, R., Narasimhan, B., Tay, K., Simon, N. & Qian, J., (2020). *glmnet: Lasso and Elastic-Net Regularized Generalized Linear Models*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/glmnet/glmnet.pdf>
- [18] Gelman, A. (2006). Multilevel (hierarchical) modeling: What it can and cannot do. *Technometrics*, 48(3), 241–251.
- [19] Gelman, A. & Hill, J. (2007). *Data analysis using regression and multilevel/hierarchical models*. Cambridge, Cambridge University Press.
- [20] Gemelli, C. N., Hayman, J., & Waller, D. (2017). Frequent whole blood donors: understanding this population and predictors of lapse. *Transfusion*, 57(1), 108–114.
- [21] Gemelli, C. N., Carver, A., Garn, A., Wright, S. T. & Davison, T. E. (2018). Evaluation of the impact of a personalized postdonation short messaging service on the retention of whole blood donors. *Transfusion*, 58(3), 701–709.
- [22] Gilbert, P., & Varadhan, R. (2019). *numDeriv: Accurate Numerical Derivatives (Version 2016.8-1.1)*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/numDeriv/index.html>
- [23] Greene, W. H. (2012). *Econometric Analysis*. (Seventh ed.): Prentice Hall.
- [24] Guo, J., Gabry, J., & Goodrich, B. (2017). *rstan: R Interface to Stan (Version 2.19.2)*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/rstan/index.html>
- [25] Hadfield, J. (2017). *MCMCglmm: MCMC Generalised Linear Mixed Models (Version 2.25)*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/MCMCglmm/index.html>
- [26] Hall, P., Johnstone, I. M., Ormerod, J. T., Wand, M. P., & Yu, J. C. F. (2020). Fast and Accurate Binary Response Mixed Model Analysis via Expectation Propagation, *Journal of the American Statistical Association*, 115:532, 1902–1916, DOI: 10.1080/01621459.2019.1665529.
- [27] Handayani, D., Notodiputro, K. A., Sadik, K., & Kurnia, A. (2017). A comparative study of approximation methods for maximum likelihood estimation in generalized linear mixed models (GLMM). *AIP Conference Proceedings*, 1827, 020033.

- [28] Hastie, T., Tibshirani, R., & Friedman, J. (2009). Shrinkage Methods. In *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. (Second ed., 61–79): Springer.
- [29] Huang, A., & Wand, M. P. (2013). Simple Marginally Noninformative Prior Distributions for Covariance Matrices. *Bayesian Analysis*, 2(2), 439–452.
- [30] Jacob, L., Obozinski, G. & Vert, J. P. (2009). Group lasso with overlap and graph lasso. 26th Annual International Conference on Machine Learning Conference Proceedings, 433–440.
- [31] Kim, A. S. I. & Wand, M. P. (2016). The Explicit Form of Expectation Propagation for a Simple Statistical Model. *Electronic Journal of Statistics*, 10(1), 550–581.
- [32] Kim, A. S. I. & Wand, M. P. (2017). On expectation propagation for generalised, linear and mixed models. *Australian and New Zealand Journal of Statistics*, 60(1), 75–102.
- [33] Kristensen, K. (2018). TMB: Template Model Builder: A General Random Effect Tool Inspired by 'ADMB' (Version 1.7.15). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/TMB/index.html>
- [34] Lauritzen, S. L. (1992). Propagation of probabilities, means and variances in mixed graphical association models. *Journal of American Statistical Association*, 87(420), 1098–1108.
- [35] Lele, S. R., Dennis, B. & Lutscher, F. (2007). Data cloning: easy maximum likelihood estimation for complex ecological models using Bayesian Markov chain Monte Carlo methods. *Ecology Letters*, 10(7), 551–563.
- [36] Liu, Q. & Pierce, D. A. (1994). A Note on Gauss-Hermite Quadrature. *Biometrika*, 81(3), 624–629.
- [37] Luts, J., & Wand, M. (2015). Variational Inference for Count Response Semiparametric Regression. *Bayesian Analysis*, 10(4), 991–1023.
- [38] Magnus, J. R. & Neudecker, H. (1999). Kronecker products, the vec operator and the Moore-Penrose inverse. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. 3 ed.: Wiley, 112.
- [39] Magnusson, A., Skaug, H., Nielsen, A., Berg, C., Kristensen, K., Maechler, M., ... Bolker, B. (2019). glmmTMB: Generalized Linear Mixed Models using Template Model Builder (Version 0.2.3). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/glmmTMB/index.html>
- [40] Spekman, M. L. C., van Tilburg, T. G. & Merz, E. M. (2019). Do deferred donors continue their donations? A large-scale register study on whole blood donor return in the Netherlands. *Transfusion*, 59(12), 3657–3665.
- [41] McCulloch, C. E., Searle, S. R., & Neuhaus, J. M. (2008). *Generalized, Linear, and Mixed Models*, New York, John Wiley & Sons.
- [42] McCulloch, C. E., & Neuhaus, J. M. (2012). Prediction of random effects in linear and generalized linear models under model misspecification. *Biometrics*, 67(1), 270–279.
- [43] Minka, T. P. (2001). Expectation propagation for approximate Bayesian inference. Seventeenth Conference on Uncertainty in Artificial Intelligence Conference Proceedings, 362–369.
- [44] Minka, T. P. 2005. Divergence measures and message passing. Microsoft Research Technical Report Series, MSR-TR-2005-173, 1–17.

- [45] Minka, T., Guiver, J., Winn, J. & Zaykov, Y. 2014. Infer.NET 2.6. Microsoft Research Cambridge. <http://research.microsoft.com/infernet>
- [46] Monahan, J. F., & Stefanski, L. A. (1989). Normal scale mixture approximations to $F^*(z)$ and computation of the logistic-normal integral. In N. Balakrishnan (Ed.), *Handbook of the Logistic Distribution*. 123, 529–540.
- [47] Nolan, T. H., & Wand, M. P. (2017). Accurate logistic variational message passing: algebraic and numerical details. *Stat*, 6(1), 102–112.
- [48] Novomestky, F. (2020). *matrixcalc*: Collection of functions for matrix calculations (Version 1.0-3). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/matrixcalc/index.html>
- [49] Ogden, H. E. (2015). A sequential reduction method for inference in generalized linear mixed models. *Electronic Journal of Statistics*, 9(1), 135–152.
- [50] Ogden, H. (2019). *glmmsr*: Fit a Generalized Linear Mixed Model (Version 0.2.3). Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/glmmsr/index.html>
- [51] Opper, M. (1999). A Bayesian Approach to On-line Learning. In D. Saad (Ed.), *On-Line Learning in Neural Networks* (363–378): Cambridge University Press.
- [52] Opper, M. & Winther, O. (2000). Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12(11), 2655–2684.
- [53] Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, California, Morgan Kaufmann Publishers, Inc.
- [54] Pebley, A., Goldman, N., & Rodríguez, G. (1996). Prenatal and delivery care and childhood immunization in Guatemala: Do family and community matter? *Demography*, 33(2), 231–247.
- [55] Pinheiro, J. C., & Bates, D. M. (2000). *Mixed-Effects Models in S and S-PLUS*, New York, Springer.
- [56] R Core Team (2019). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- [57] Raudenbush, S. W. (1993). A Crossed Random Effects Model for Unbalanced Data with Applications in Cross-Sectional and Longitudinal Research. *Journal of Educational Statistics*. 18(4), 321–349.
- [58] Rohde, D. & Wand, M. P. (2015). Semiparametric Mean Field Variational Bayes: General Principles and Numerical Issues. *Journal of Machine Learning Research*, 17(172), 1–47.
- [59] Stan Development Team (2020). *RStan: the R interface to Stan*. R package version 2.21.2. <http://mc-stan.org/>.
- [60] Scott, M. A., Simonoff, J. S., & Marx, B. D. (2013). *The SAGE Handbook of Multilevel Modeling*.
- [61] Solymos, P. (2010). *dclone*: Data Cloning in R. *The R Journal*, 2(2), 29–37.
- [62] Steenbergen, M. R., & Jones, B. S. (2002). Modeling Multilevel Data Structures. *American Journal of Political Science*, 46(1), 218–237.

-
- [63] Teunissen, P. J. G. (2007). Best prediction in linear models with mixed integer/real unknowns: theory and application. *Journal of Geodesy* 81, 759–780.
- [64] Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1), 267–288.
- [65] Vonesh, E. F. (1996). A note on the use of Laplace’s approximation for nonlinear mixed-effects models. *Biometrika*, 83(2), 447–452.
- [66] Wainwright, M. J. & Jordan, M. I. (2008). Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends in Machine Learning*, 1.
- [67] Wand, M. P & Jones, M. C. (1994). *Kernel Smoothing*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability (60). Boca Raton, Chapman & Hall.
- [68] Wand, M. P. (2017). Fast Approximate Inference for Arbitrarily Large Semiparametric Regression Models via Message Passing. *Journal of the American Statistical Association*, 112(517), 137–168.
- [69] Wand, M. P. & Yu, J. C. F. (2020). *glmmEP: Generalized Linear Mixed Model Analysis via Expectation Propagation (Version 1.0-3.1)*. Comprehensive R Archive Network. Retrieved from <https://cran.r-project.org/web/packages/glmmEP/index.html>
- [70] Winn, J. & Bishop, C. (2005). Variational Message Passing. *Journal of Machine Learning Research*, 6, 661–694
- [71] Wolfinger, R., & O’connell, M. (1993). Generalized linear mixed models a pseudo-likelihood approach. *Journal of Statistical Computation and Simulation*, 48, 233–243.