

# FAST APPROXIMATE INFERENCE FOR LONGITUDINAL AND MULTILEVEL DATA ANALYSIS

Cathy Yuen Yi Lee

GStat MBiostat BMathAdv(Hons)

*A PhD thesis submitted in fulfilment of the requirements for  
Doctor of Philosophy in Mathematics  
in the School of Mathematical and Physical Sciences,  
University of Technology Sydney*

2016

© Cathy Yuen Yi Lee, 2016. All rights reserved.

Permission is herewith granted to University of Technology Sydney to circulate and to have copied for non-commercial purposes, at its discretion, the above title upon request of individuals and institutions.

---

## Certificate of Original Authorship

I, Cathy Yuen Yi Lee, certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of student:

Date:

---

## Acknowledgements

*Fear protects us in many ways, but what has served me is diving into my own obliviousness. Being more confident than I should be trying things that you never would have tried. Your inexperience is an asset in that it will make you think in original, unconventional ways. Accept your lack of knowledge and use it as your asset.*

Natalie Portman

I never thought I would go to university, let alone complete a PhD degree. I am always and forever grateful to my lovely parents for instilling the value of education in me and for making countless sacrifice to always provide for their children. I wouldn't be who I am today if it wasn't for the undying love and support my parents have unconditionally given to me.

Words are simply not enough to express my sincere gratitude to my academic parents, Matt and Louise, for believing in me and seeing my good qualities way before I ever could. Through watching them, I learned dedication and perseverance to succeed. I am no longer afraid of following my dream I had set since I was 16. Thank you for guiding me to reach a stage where I can be proud of who I am and what I have achieved for the past three years.

I must express my gratitude to Jason, my companion, for his enthusiasm and continuous support and encouragement. I was amazed by his willingness and determination to support me in every possible way throughout this journey. He continually challenged me to dig deeper and learn more, and I greatly appreciate everything he has done to make this possible.

“Friends are the bold essence of life”. Completing this work would have been more difficult were it not for the support and friendship provided by my colleagues at the School of Mathematical and Physical Sciences. A special thanks to Chris, who has proofread countless pages of my work. Also a special thanks to Marianne, who has always been there for me and experienced all of the laughter and tears from my special journey.

I finish here with Hon's favourite quote “Stay Hungry, Stay Foolish”. I believe this quote will continue to motivate me in life for many years to come.

---

## Publications

Lee, C. Y. Y. and Wand, M. P. (2016). Streamlined mean field variational Bayes for longitudinal and multilevel data analysis. *Biometrical Journal*. DOI: 10.1002/bimj.201500007

Bentley, J. and Lee, C. Y. Y. (2016). Bayesian Inference for Gaussian Semiparametric Multilevel Models. *SAS Global Forum*, April 18-21, Las Vegas, United States. <http://support.sas.com/resources/papers/proceedings16/7200-2016.pdf>

Lee, C. Y. Y. and Wand, M. P. (2015). Variational methods for fitting complex Bayesian mixed effects models to health data. *Statistics in Medicine*. DOI: 10.1002/sim.6737.

Lee, C. Y. Y., Homer, C., Bisits, A. and Ryan, L. (2015). Increasing variation in hospital caesarean section rates among low-risk nulliparous women in Australia, from 1994 to 2010. *Birth in revision*.

Chen, J., Lee, C. Y. Y., Coull, B., Valeri, L., Gruskin, S., Beckfield, J., Singh, N. and Krieger, N. (2015). 50-year trends and state variation in socioeconomic and racial/ethnic inequalities in US infant death rates, 1960 to 2010. *Under preparation*.

---

The last listed publication related to my 3-week academic visit at the Harvard School of Public Health during March-April 2014. Throughout my visit, I collaborated with a research team of world-renowned social epidemiologists, biostatisticians and research scientists on a study under an overarching project entitled *The Public Health Disparities Geocoding Project Monograph*. For more information, please refer to <https://www.hsph.harvard.edu/thegeocodingproject/>.

---

## Presentations

Lee, C. Y. Y and Bentley, J. P. Fast variational Bayesian inference for Gaussian semi-parametric multilevel models in SAS/IML<sup>®</sup>, Las Vegas, United States, 2016. (Accepted for a 20-min workshop presentation)

Lee, C. Y.Y. Fast approximate inference for longitudinal and multilevel data analysis. *The 16th J.B. Douglas Awards, Statistical Society of Australia*, Sydney, Australia, 2015. (Awarded Equal First Price for oral presentation)

Lee, C. Y. Y., Wand, M. P. and Ryan, L. Fast approximate inference for longitudinal and multilevel data analysis. *The 16th J.B. Douglas Awards, Statistical Society of Australia*, Sydney, Australia, 2015. (Awarded Equal First Price for oral presentation)

Lee, C. Y. Y., Chen, J., Coull, B., Valeri, L., Gruskin, S., Beckfield, J., Singh, N. and Krieger, N. 50-year trends and state variation in socioeconomic and racial/ethnic inequalities in US infant death rates, 1960 to 2010. *Joint Statistical Meeting*, Seattle, United States, 2015.

Lee, C. Y. Y. and Wand, M. P. Fast approximate variational inference for Bayesian multilevel models. *Young Statisticians Conference*, Adelaide, Australia, 2015. (Awarded 3rd Price for oral presentation)

Ryan, L. and Lee, C. Y. Y. Spatio-temporal variation in cesarean rates in New South Wales, Australia, from 1994 to 2010. *Joint Statistical Meeting*, Boston, United States, 2014.

Lee, C. Y. Y. and Wand, M. P. Fast approximate variational inference for Bayesian multilevel models. *Australian Statistical Conference in conjunction with the Institute of Mathematical Statistics Annual Meeting*, Sydney, Australia, 2014.

---

The fourth listed presentation related to my 3-week academic visit at the Harvard School of Public Health during March-April 2014. Throughout my visit, I collaborated with a research team of world-renowned social epidemiologists, biostatisticians and research scientists on a study under an overarching project entitled *The Public Health Disparities Geocoding Project Monograph*. For more information, please refer to <https://www.hsph.harvard.edu/thegeocodingproject/>.

---

Lee, C. Y. Y. and Wand, M. P. Fast approximate variational inference for Bayesian longitudinal and multilevel semiparametric regression models. *University of Queensland*, Brisbane, Australia, 2014.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Aim and motivation . . . . .	1
1.2	Structure of the thesis . . . . .	3
1.3	Semiparametric regression . . . . .	4
1.3.1	Mixed model representation . . . . .	6
1.4	Graphical models . . . . .	6
1.5	Bayesian inference . . . . .	9
1.5.1	A Bayesian viewpoint of probability and statistics . . . . .	9
1.5.2	Markov chain Monte Carlo . . . . .	10
1.5.3	A brief introduction to variational approximations . . . . .	11
1.6	Terminological and notational conventions . . . . .	17
1.7	Common distributions . . . . .	18
1.8	Distributional results . . . . .	18
1.9	Vector differential calculus . . . . .	20
1.10	Special functions . . . . .	20
1.11	Matrix results . . . . .	21
1.12	Statistical software . . . . .	22
<b>2</b>	<b>Mean Field Variational Bayes Approximations for Longitudinal and Multilevel Data Analysis</b>	<b>23</b>
2.1	Introduction . . . . .	23
2.2	A brief introduction to longitudinal and multilevel models . . . . .	25
2.2.1	Graph theoretical viewpoint of longitudinal and multilevel models . . . . .	27
2.2.2	Random intercepts models . . . . .	27
2.2.3	Random intercepts and slopes models . . . . .	28
2.2.4	Extension to semiparametric regression . . . . .	29
2.3	Mixed model representation . . . . .	29
2.3.1	Sample size and subscript notation . . . . .	31
2.4	Prior distributions for fixed and random effects . . . . .	33

2.5	Gaussian semiparametric mixed models . . . . .	34
2.5.1	Approximate Bayesian inference via mean field variational Bayes . .	34
2.5.2	Approximating density functions of entries of $\Sigma^R$ . . . . .	39
2.6	Student- $t$ semiparametric mixed models . . . . .	41
2.6.1	Approximate Bayesian inference via mean field variational Bayes . .	41
2.7	Bernoulli semiparametric mixed models . . . . .	45
2.7.1	Approximate Bayesian inference via mean field variational Bayes . .	46
2.8	Poisson semiparametric mixed models . . . . .	49
2.8.1	Approximate Bayesian inference via non-conjugate variational mes- sage passing . . . . .	49
2.9	Displaying approximate posterior means of regression function fits . . . . .	55
2.10	Numerical evaluation . . . . .	56
2.10.1	Fitting a Bayesian hierarchical model in Stan . . . . .	57
2.10.2	Assessment of accuracy . . . . .	59
2.10.3	Assessment of coverage . . . . .	60
2.10.4	Assessment of speed . . . . .	60
2.11	Concluding remarks . . . . .	64
2.A	Optimal $q$ -densities derivation for Gaussian semiparametric mixed models .	65
2.A.1	Derivation of the marginal log-likelihood lower bound . . . . .	71
2.B	Optimal $q$ -densities derivation for Student- $t$ semiparametric mixed models .	75
2.B.1	Derivation of the marginal log-likelihood lower bound . . . . .	76
2.C	Optimal $q$ -densities derivation for Bernoulli semiparametric mixed models .	80
2.C.1	Derivation of the marginal log-likelihood lower bound . . . . .	81
2.D	Optimal $q$ -densities derivation for Poisson semiparametric mixed models . .	82
2.D.1	Derivation of the marginal log-likelihood lower bound . . . . .	84
<b>3</b>	<b>Streamlining Mean Field Variational Bayes Algorithms And Three-Level Model Extensions</b>	<b>85</b>
3.1	Introduction . . . . .	85
3.2	Computational challenges in naïve MFVB algorithms . . . . .	86
3.3	Predictor structure and matrix notation . . . . .	88
3.4	Streamlining mean field variational Bayes algorithms . . . . .	89
3.4.1	Streamlining update expressions involving $\Sigma_{q(\beta, \mathbf{u})}$ . . . . .	91
3.4.2	Streamlining lower bound expression involving $\Sigma_{q(\beta, \mathbf{u})}$ . . . . .	93
3.5	Computational speed: naïve versus streamlined . . . . .	94
3.6	Real data applications . . . . .	101
3.6.1	Application to smoking data . . . . .	101
3.6.2	Application to student assessment data . . . . .	103
3.6.3	Application to German health care data . . . . .	104



3.7	Extension to three-level semiparametric mixed models . . . . .	105
3.7.1	Mixed model representation . . . . .	110
3.7.2	Naïve mean field variational Bayes algorithms . . . . .	113
3.7.3	Streamlined mean field variational Bayes algorithms . . . . .	114
3.8	Concluding remarks . . . . .	114
<b>4</b>	<b>Measurement Error, Missing Data and Real-Time Extensions</b>	<b>118</b>
4.1	Introduction . . . . .	118
4.2	Gaussian semiparametric mixed models with measurement error problems .	119
4.2.1	Notation . . . . .	121
4.2.2	Approximate Bayesian inference via mean field variational Bayes . .	124
4.3	Gaussian semiparametric mixed models with missing data problems . . . .	126
4.3.1	Notation . . . . .	129
4.3.2	Approximate Bayesian inference via mean field variational Bayes . .	130
4.4	Numerical evaluation . . . . .	134
4.4.1	Assessment of accuracy . . . . .	135
4.4.2	Assessment of coverage . . . . .	136
4.4.3	Assessment of speed . . . . .	136
4.5	Real-time mean field variational Bayes algorithms . . . . .	136
4.6	Concluding remarks . . . . .	142
4.A	Optimal $q$ -densities derivation for measurement error problems . . . . .	147
4.A.1	Derivation of the marginal log-likelihood lower bound . . . . .	154
4.B	Optimal $q$ -densities derivation for missing data problems . . . . .	157
<b>5</b>	<b>Extension to Group-Specific Curve Models With Contrasting</b>	<b>161</b>
5.1	Introduction . . . . .	161
5.2	Mean field variational Bayes approximations to caesarean section data . . .	164
5.2.1	Random intercept and slope model . . . . .	166
5.2.2	Group-specific curve model . . . . .	170
5.2.3	Factor-by-curve interactions . . . . .	171
5.3	Numerical evaluation . . . . .	176
5.3.1	Assessment of accuracy . . . . .	176
5.3.2	Assessment of speed . . . . .	177
5.4	Application to caesarean section data . . . . .	178
5.5	Concluding remarks . . . . .	182
5.A	Derivation of the mean field variational Bayes algorithm for model (5.11) .	185
5.A.1	Mixed model representation . . . . .	186
5.A.2	Streamlined mean field variational Bayes algorithm . . . . .	188
5.A.3	Fitting the group-specific curve model (5.11) in <code>Rstan</code> . . . . .	188

<b>6</b>	<b>Alternative Approach Based on Variational Message Passing</b>	<b>193</b>
6.1	Introduction . . . . .	193
6.2	Natural canonical forms for exponential family densities . . . . .	194
6.3	Primitive integrals and function definitions . . . . .	195
6.4	Factor graphs . . . . .	197
6.5	Variational message passing . . . . .	198
6.6	Illustrative example of simple mixed effects regression . . . . .	200
6.7	Arbitrarily large models viewpoint . . . . .	205
6.8	Concluding remarks . . . . .	209
<b>7</b>	<b>Conclusions and Future Directions</b>	<b>211</b>

# List of Figures

1.1	Directed acyclic graph involving six random variables: $x_1, \dots, x_6$ . The open circles represent the random (“hidden”) nodes. The shaded circle represents the observed (“evidence”) node. The grey dashed line polygon represents the Markov blanket of the node $x_4$ . . . . .	8
1.2	Left: The smallest ancestral subgraph corresponding to Figure 1.1. Right: The moral graph corresponding to Figure 1.1. . . . .	8
1.3	Factor graph corresponding to the DAG in Figure 1.1. . . . .	15
2.1	Panel of scatterplots of the 11-year old mathematics test scores versus the eight-year old test scores. Each panel corresponds to a school. Details of the data are described in Goldstein (2010). . . . .	26
2.2	Directed acyclic graph corresponding to model (2.1). The open circles represent the random nodes. The shaded circle represents the observed node. The vector $\mathbf{y} = (y_{11}, \dots, y_{mn_m})$ is represented by a single node, and $\beta_0$ and $\beta_x$ are represented by the $\beta$ node. Similarly, $\mathbf{u}^R = (u_1^R, \dots, u_m^R)$ . . . . .	28
2.3	Directed acyclic graph for the two-level Bayesian semiparametric mixed models with the Gaussian and Student- $t$ responses. The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The pink open nodes $\mathbf{b}$ and $\nu$ are additional parameters for the Student- $t$ response model. The numbered grey dashed polygons indicate the corresponding update expressions in Algorithm 3. . . . .	35
2.4	Directed acyclic graph for the two-level Bayesian semiparametric mixed models with the Bernoulli and Poisson responses. The shaded node corresponds to the observed data vector and the open node corresponds to the random or auxiliary variables. The numbered grey dashed polygons indicate the corresponding update expressions in Algorithms 4 and 6. . . . .	35

2.5	Life cycle of nodes for mean field variational Bayes: (a) directed acyclic graph for model (2.11), (b) moral graph for model (2.11), (c) modification of Figure 2.5b with eight edges removed to impose $q$ -density product restriction. In each graph, shading is used to signify the observed data vector.	37
2.6	Variational representation of the function $-\log(1 + e^x)$ as the maximum of a family of parabolas, corresponding to (2.20).	47
2.7	Summary of MCMC samples for fitting the two-level Bayesian semiparametric mixed model with Gaussian response to the simulated data. The columns are: parameter name, trace plot of the MCMC sample, plot of sample against its lag-one sample, sample autocorrelation function, kernel density estimate of posterior density function and numerical summaries of posterior density function.	61
2.8	Approximate posterior density functions for the covariance parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication from each of the simulation studies described in the text. The green vertical lines represent the true parameter values. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against a MCMC benchmark.	62
2.9	Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication of the simulation study described in the text. The red curves represent the true mean functions and the sky blue circles represent the simulated data.	63
2.10	Side-by-side boxplots of accuracy scores for MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed models.	63
3.1	Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the two-level Bayesian semiparametric mixed models to the real datasets. The sky blue circles represent the real data.	106
3.2	Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Gaussian response model to the smoking data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark.	107

3.3	Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Bernoulli response model to the Program of International Student Assessment data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark.	108
3.4	Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Poisson response model to the German health care data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark. . . . .	109
3.5	Directed acyclic graphs for model (3.17). The shaded node corresponds to the observed data vector. The colour keys at the top of the figure indicate the components of the graph corresponding to each response outcome. . . .	112
3.6	Different $\mathbf{Z}$ structures for the two-level and three-level models with $m = 2$ . The $\mathbf{Z}$ matrix starts off as having a simple block-diagonal structure and, as the level of hierarchy increases, it grows into a “nested” block-diagonal structure. The definitions of matrices are described in the text. . . . .	113
4.1	Directed acyclic graph corresponds to model (4.1). The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are presented in the Gaussian model DAG as shown in Figure 2.3. . . . .	122
4.2	Directed acyclic graph corresponds to model (4.7). The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are presented in the Gaussian model DAG as shown in Figure 2.3. . . . .	122
4.3	Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed model with Gaussian response subject to classical measurement errors. . . . .	137
4.4	Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed model with Gaussian response subject to missingness. . . . .	138
4.5	Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication of the simulation study described in the text. The red curves represent the true mean functions. The grey circles represent the observed data and the sky blue circles represent the unobserved or missing data. . . . .	139

4.6 Diagrammatic description of the real-time/online variational analysis. The first phase is batch phase, where we start by running a small subset of the initial data in the batch algorithm to initialise model parameters and obtain starting values for the data sufficient statistics (these are all we keep). The second phase is the real-time/online phase where the MFVB algorithm updates as each new data point arrives. Real-time updates use only the new data and summary statistics from previous iterations rather than the full set of available data. . . . . 141

4.7 Real-time MFVB fitting of the two-level Bayesian semiparametric mixed models with Gaussian response for two different sample sizes. The sky blue circles represent the batch data and the pink circles represent the new data. 145

5.1 Trends in caesarean section rates for low-risk nulliparous women in the largest state of Australia, New South Wales, between 1994 and 2010. The dotted line is the overall mean trend and the solid lines are the selected hospital-specific trends. . . . . 162

5.2 Trends in caesarean section rates for low-risk nulliparous women aged less than 25 years and those aged greater or equal to 25 years in New South Wales, Australia, from 1994 to 2010. The dotted line is the overall mean trend and the solid lines are the selected hospital-specific trends. . . . . 172

5.3 Various structures of the  $\mathbf{Z}$  matrix across logistic mixed models (5.4), (5.9) and (5.11), with the number of groups  $m = 2$ . The  $\mathbf{Z}$  matrix starts off as having a simple block-diagonal structure and, as the model increases in complexity, it grows into a “nested” block-diagonal structure. The definitions of matrices are described in Subsection 5.2.2. . . . . 175

5.4 The approximate posterior density functions obtained from MFVB and MCMC for a single replication of the simulation study described in the text. Each pair of density function corresponds to a model parameter  $f(Q_k)$ ,  $1 \leq k \leq 3$  and  $g_i(Q_k)$ ,  $k = 2$ , where  $Q_k$  is the  $k$ th sample quintile of the  $x$ s. The vertical lines represent the true parameter values. The accuracy scores on the top right of on each plot show the accuracy of MFVB approximation compared against a MCMC benchmark. . . . . 178

5.5 Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs. Each boxplot corresponds to a model parameter  $f(Q_k)$ ,  $1 \leq k \leq 3$  and  $g_i(Q_k)$ ,  $k = 2$ , where  $Q_k$  is the  $k$ th sample quintile of the  $x$ s. . . . . 179

5.6	The MFVB fitted hospital-specific probability functions of caesarean section for low-risk nulliparous women of aged less than 25 years and those of aged greater or equal to 25 years, as a function of time for each hospital. The dashed curves represent pointwise 95% credible sets. Each panel corresponds to a different hospital. . . . .	180
5.7	The MFVB estimated overall and selected hospital-specific contrast curves defined by (5.12), corresponding to the odds of caesarean section for low-risk nulliparous women aged greater or equal to 25 years, as a function of time, compared with those aged less than 25 years. The shaded regions correspond to pointwise 95% credible sets. . . . .	181
5.8	The MFVB estimated hospital-specific odds ratios of caesarean section for low-risk nulliparous women of older age compared with those of younger age in the year 2010. . . . .	182
5.9	Directed Acyclic Graph for the Bernoulli group-specific curve model with contrasting. The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are present in the Bernoulli model DAG graph as shown in Figure 2.4. . . . .	186
6.1	A factor graph corresponding to the function $g(x_1, x_2, x_3, x_4, x_5)$ defined by (6.3). . . . .	197
6.2	A factor graph corresponding to a Bayesian model with stochastic nodes $\theta_1, \dots, \theta_9$ and factors $f_1, \dots, f_{10}$ . . . . .	199
6.3	Each of the messages between neighbouring nodes on the factor graph for the random intercepts and slopes model illustrative example. The pink arrows depict the direction from the stochastic node to factor, while the grey arrows depict the direction from the factor to stochastic node. . . . .	201
6.4	Diagrammatic depiction of the extension from simple mixed effects regression to three-level penalised spline mixed regression. The fragments shown in darker blue are the ones that are additional to the fragments in the simple mixed model. The fragments shown in grey are presented in Figure 6.3. . . . .	208

# List of Tables

1.1	Distributions used in this thesis and their corresponding probability density or mass functions. . . . .	19
2.1	Average (standard error) elapsed of the computing times in seconds for MCMC and MFVB fitting of the two-level Bayesian semiparametric mixed models to the simulated data. . . . .	60
2.2	Percentage coverage of true parameter values by approximate 95% credible sets based on the MFVB approximate posterior density functions. The percentages are based on 1000 replications. . . . .	64
3.1	Average (standard error) elapsed of the computing times in seconds for the simulation described in the text, using the naïve Algorithm 3, streamlined Algorithm 7 and <code>gam()</code> function in R <code>mgcv</code> package. The ratio of naïve over streamlined and the ratio of <code>gamm()</code> over streamlined are also presented. . .	101
3.2	Description of the United States National Center for Health Statistics perinatal health data as presented in Abrevaya (2006). . . . .	102
3.3	Description of the 2000 Program for International Student Assessment conducted by the Organisation for Economic Cooperation and Development. .	103
3.4	Description of the German health care data as presented in Winkelmann (2004). . . . .	104
4.1	Percentage coverage of true parameter values by approximate 95% credible sets based on MFVB approximate posterior density functions. A value of $RR = 0.9$ corresponds to a small amount of measurement error and $RR = 0.7$ corresponds to a substantial corruption of the predictor. Low missingness for the MNAR model corresponds to $(\phi_0, \phi_1) = (2.95, -2.95)$ and high missingness to $(\phi_0, \phi_1) = (0.85, -1.05)$ . The percentages are based on 1000 replications. . . . .	140



*LIST OF TABLES*

---

4.2	Average (standard error) elapsed of the computing times in seconds for the MFVB and MCMC fitting of the two-level Bayesian semiparametric mixed models with Gaussian response subject to classical measurement error or missing data problem. . . . .	140
5.1	Average (standard error) elapsed of the computing times in seconds for the streamlined MFVB Algorithm15 in the simulation setting described in (5.14).177	
5.2	The MFVB estimated odds ratios of caesarean section (averaged across hospitals) for low-risk nulliparous women aged greater than or equal to 25 years compared with those aged less than 25 years for selected years of birth.182	
6.1	Expressions for natural statistics, natural parameter vectors and natural canonical forms of common exponential family densities. . . . .	195

# Abstract

Generalised linear mixed models are the cornerstone of longitudinal and multilevel data analysis. However, exact inference for Bayesian mixed models with semiparametric extensions is typically intractable, requiring approximate inference methods for use in practice. Markov chain Monte Carlo or MCMC is one of the most commonly used approximate inference methods in this setting, but can be computationally intensive and often suffers from poor convergence in complex models. A faster, deterministic alternative to MCMC is variational approximations, a class of deterministic algorithms that is based on reformulating the problem of computing the posterior distribution as an optimisation problem, simplifying that problem and finding solutions to the perturbed problem. In this thesis, we work with a particular class of variational approximations, known as the mean field variational Bayes (MFVB). In essence, MFVB approximations are based upon optimising the Kullback-Leibler divergence with respect to the so-called approximating distribution. We derive MFVB algorithms for a wide variety of Bayesian semiparametric mixed models with Gaussian, Student- $t$ , Bernoulli and Poisson responses. In order to overcome the computational cost of the direct naïve approach to the underlying MFVB calculations for models, we introduce a novel, streamlined approach that involves matrix permutation and block decomposition. Through a series of numerical studies, we demonstrate that the MFVB algorithms achieve a good level of accuracy compared to a MCMC benchmark (our gold standard). Furthermore, our developed streamlined algorithms are shown to have a complexity that is linear in the number of groups at each level, representing a two orders of magnitude improvement over the naïve approach. More importantly, the modularity of MFVB allows relatively simple extensions to more complicated scenarios, including higher-level random effects, measurement error and/or missing data problems, models with group-specific curves and real-time or online data processing. Illustrations from various real data examples are provided.

# Chapter 1

## Introduction

*The world is being re-shaped by the convergence of social, mobile, cloud, big data, community and other powerful forces. The combination of these technologies unlocks an incredible opportunity to connect everything together in a new way and is dramatically transforming the way we live and work.*

Marc Benioff

### 1.1 Aim and motivation

In today's digital age, we are surrounded by a vast amount of data, so called "Big Data". In a recent talk at Google, Kenneth Cukier, Data Editor of *The Economist* stated that Big Data is a revolution that will transform how we live, work and think about the economy, health and society in the years to come. "The amount of data in the world is increasing exponentially, and we don't necessarily have the tools to handle it.", Kenneth says. So the question is, how can we take this opportunity and turn this new wealth of information into big insights? The answer is *fast data analysis*. The key role of data analysis is to translate the raw, intimidating numbers into plain english, scientific findings. "The doubling of computing power every 18 months (Moore's Law) is nothing compared to a big algorithm - a set of rules that can be used to solve a problem a thousand times faster than conventional computational methods could", says Albert J. Weatherhead III University Professor Gary King at Harvard University. Inspired by these people, I am motivated to take on the Big Data challenges that show a presence within this thesis.

Most datasets nowadays in science have some form of grouped or hierarchical structure. Such structure may arise from repeated measurements of individuals over time, students grouped within classrooms and schools, or friends grouped within networks (much like LinkedIn or Facebook). The reality of scientific research is that our fundamental units of interest are usually intertwined, and quite often, of primary interest from a theoretical

perspective. Insightful information can be gleaned from connecting gene expression data, linked administrative health records or environmental sequences, but how can we best extract signals from these correlated data? From a methodological perspective, the complex structure of these data complicates the conventional statistical methods and threatens our ability to fitting models and performing statistical inference efficiently.

Bayesian generalised linear mixed models (e.g. Gelman and Hill, 2007) continue to grow in popularity in response to the demands of large grouped/clustered datasets and have proved powerful in many real-world applications. These models are based on the deceptively simple Bayes's Theorem. However, the computations often involve intractable multivariate integrals in likelihood and posterior density expressions, requiring approximate inference methods for use in practice. *Markov chain Monte Carlo* or MCMC is one of the most prominent methods for estimating Bayesian statistical models and has proven useful in a wide array of problems. Unfortunately, in many situations a MCMC type of sampling strategy is undesirably slow and suffers from poor mixing. In parallel to these challenges in statistics, the computer science community has been developing approximate solutions to inferential problems using variational algorithms (e.g. Bishop, 2006; Ormerod and Wand, 2010). These are deterministic algorithms that facilitate analytical calculations for the posterior distributions, an iterative approach resembling an Expectation-Maximisation method. In contrast to these algorithms, MCMC methods yield their results in the form of a set of samples from the posteriors. Practical implementations of variational algorithms make use of factorised approximating posteriors and priors that belong to the conjugate-exponential family, making the required integrals tractable. The resultant variational approximations therefore sacrifice some accuracy of MCMC, albeit minimal, but offer vast improvements in terms of speed and memory efficiency.

What is currently lacking in the statistical literature is research on the accuracy of the variational approximations for Bayesian generalised linear mixed models with various extensions. All of these bring the overarching question of my thesis:

*Are variational approximations a viable set of tools for Bayesian longitudinal and multilevel data analysis?*

To address this we: first, identify models for which standard inference is impaired by time constraints; second, develop variational algorithms; and third, empirically assess the performance of variational algorithms in terms of inferential accuracy and computational speed. In what follows, we commence our discussion by presenting the literature review on approximate Bayesian inference.

## 1.2 Structure of the thesis

Each chapter in this thesis is dedicated to the development of variational algorithms catered to a particular longitudinal or multilevel model setting. In Chapter 2, we consider the generalised linear mixed models for longitudinal and multilevel data analysis. Four common response distributions, including Gaussian, Student- $t$ , Bernoulli and Poisson, are explored. We provide a step-by-step, notationally friendly guide on the development of fast variational algorithms for fitting and inference in mixed models, where all algorithmic updates are available in closed form. The advantages of variational methods as compared to the gold standard MCMC methods in terms of speed, coverage and accuracy are illustrated via a series of comprehensive numerical studies. A computational challenge in the directly-implemented variational algorithms arises from the update expression of the effects covariance matrix, which involves inversion and storage of a large sparse matrix, a cost that grows cubically in the number of groups. The centrepieces for Chapter 3 are our developed novel, streamlined variational algorithms for fast and memory efficient fitting and inference in large Bayesian generalised linear mixed models with semiparametric extensions - currently represent the state-of-the-art algorithms in this area. The number of operations is linear in the number of groups at each level, which constitutes a major improvement over the naïve/direct approach. Storage requirements are also lessened considerably.

The beauty of the Bayesian paradigm combined with variational algorithms is its tremendous flexibility. In Chapter 4, we extend our variational algorithms to cater for three interesting and challenging scenarios, including data subject to measurement error problems, data subject to missing data problems, and data that are processed in real-time as they arrive. Accuracy and time comparisons are considered as an effective assessment of the developed algorithms. Chapter 5 is motivated by a study that examines trends in caesarean section rates in the largest state of Australia, New South Wales, between 1994 and 2010. We propose a group-specific curve model that encapsulates the complex non-linear features of the overall and hospital-specific trends in caesarean section rates, while taking into account hospital variability over time. We use penalised spline based smooth functions that represent trends and implement a fully variational Bayes approach to model fitting. In Chapter 6, we present a discussion on a variant of the standard variational Bayes approach that is based on natural parameters and sufficient statistics and allows an easier extension to arbitrarily large models via factor graph fragments. We demonstrate this approach for the two models considered in Chapters 2 and 3.

The majority of the materials written in this thesis have been published (one is under review at the time of writing) and presented at five conferences. Main results in Chapters 2, 3 and 5 have been published in Lee and Wand (2015a), Lee *et al.* (2015) and Lee and Wand (2015b) respectively.

### 1.3 Semiparametric regression

Parametric regression methods for longitudinal and multilevel data analysis have been well developed in the past two decades. Such methods can be broadly classified into estimating equations based methods (e.g. Breslow and Clayton, 1993; Zeger and Liang, 1986) and mixed models (e.g. Gelman and Hill, 2007; Goldstein, 2011; Hastie and Tibshirani, 1990). Parametric regression assumes a functional form (typically linear) in the relationship between the mean of a response and predictors. Although such parametric assumption offers simplicity, it is inappropriate for situations when the relationship between the mean response and predictors is unknown.

Semiparametric regression is a seamless fusion between parametric and nonparametric regression analysis (Ruppert *et al.*, 2009). It is a prominent field synthesising research across several branches of statistics including parametric and non-parametric regression, longitudinal and multilevel data analysis, and Bayesian hierarchical models. Early comprehensive reviews on semiparametric regression analysis are Ruppert *et al.* (2003) and Ruppert *et al.* (2009). Semiparametric regression models extends classical generalised linear mixed models to accomodate non-linear predictor effects flexibly using penalised basis functions such as B-splines and Daubechies wavelets. Such penalisation can be achieved through fitting random effects that have the same formulation as those used traditionally in longitudinal and multilevel data analysis. Throughout this thesis, we work with a direct generalisation of smoothing splines, namely the *O'Sullivan splines* (O'Sullivan, 1986, Section 3). O'Sullivan splines are a class of penalised splines based on B-spline basis functions. They possess the attractive feature of requiring considerably fewer basis functions and their smoothness, numerical stability and natural boundary conditions make them the most widely used class of penalised splines in standard statistical software, for example, the `smooth.spline()` function in R (R Development Core Team, 2015) and the MIXED procedure in SAS (SAS Institute Inc., 2013). Wand and Ormerod (2008) give a detailed description of O'Sullivan penalised splines and their mixed model representation. A brief sketch of their description is provided here for reference.

Consider the simplest non-parametric regression setting

$$y_i = f(x_i) + \varepsilon_i, \quad 1 \leq i \leq n,$$

where  $(x_i, y_i) \in \mathbb{R} \times \mathbb{R}$ . Suppose we are interested in estimating  $f$  over the interval  $[a, b]$  containing the  $x_i$ s via a set of cubic B-spline basis functions  $\mathbf{B}_x \equiv [B_1(x), \dots, B_{K+4}(x)]$  for  $K \leq n$ . The corresponding knot sequence is defined by  $a = \kappa_1 = \kappa_2 = \kappa_3 = \kappa_4 < \kappa_5 < \dots < \kappa_{K+4} < \kappa_{K+5} = \kappa_{K+6} = \kappa_{K+7} = \kappa_{K+8} = b$  (Hastie *et al.*, 2005). The coefficients may be estimated in a number of ways. The simplest is the *penalised residual sum of squares*

### 1.3. SEMIPARAMETRIC REGRESSION

---

or PRSS, which involves choosing the coefficients to minimise

$$PRSS(f, \lambda) = \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda \int_a^b f''(x) dx. \quad (1.1)$$

The expression  $\lambda \int_a^b f''(x) dx$  is the so-called *penalty term* because it penalises fits that are too rough, thus yielding a smoother result. The amount of smoothing is controlled by  $\lambda > 0$ , and is usually referred to as *a smoothing parameter*. The solution to (1.1) is the O'Sullivan penalised spline  $f(\mathbf{x}) = \mathbf{B}\boldsymbol{\nu}$  and thus (1.1) can be rewritten as

$$PRSS(\boldsymbol{\nu}, \lambda) = (\mathbf{y} - \mathbf{B}\boldsymbol{\nu})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\nu}) + \lambda \boldsymbol{\nu}^\top \boldsymbol{\Omega} \boldsymbol{\nu}, \quad (1.2)$$

where  $B_{ik} = B_k(x_i)$  and  $\Omega_{kk'} = \int_a^b B_k''(x) B_{k'}''(x) dx$ . Straightforward algebraic manipulation leads to the following fitted O'Sullivan penalised spline with a solution to (1.2)

$$\hat{f}(\mathbf{x}) = \mathbf{B}\hat{\boldsymbol{\nu}}, \quad \hat{\boldsymbol{\nu}} = (\mathbf{B}^\top \mathbf{B} + \lambda \boldsymbol{\Omega})^{-1} \mathbf{B}^\top \mathbf{y}. \quad (1.3)$$

Computation of the design matrix  $\mathbf{B}$  is straightforward and is readily available in the R environment. However, computation of the penalty matrix  $\boldsymbol{\Omega}$  requires some effort. In Section 6 of Wand and Ormerod (2008), an exact matrix algebraic expression for  $\boldsymbol{\Omega}$  is derived by applying the Simpson's rule over each of the inter-knot differences, given as

$$\boldsymbol{\Omega} = (\tilde{\mathbf{B}}'')^\top \text{diag}(\mathbf{w}) \tilde{\mathbf{B}}'',$$

where  $\tilde{\mathbf{B}}''$  is the  $3(K+7) \times (K+4)$  matrix with the  $(i, j)$ th entry  $\mathbf{B}_j''(\tilde{x}_i)$  and  $\tilde{x}_i$  is the  $i$ th entry of the vector

$$\tilde{\mathbf{x}} = \left( \kappa_1, \frac{\kappa_1 + \kappa_2}{2}, \kappa_2, \kappa_2, \frac{\kappa_2 + \kappa_3}{2}, \kappa_3, \dots, \kappa_{K+7}, \frac{\kappa_{K+7} + \kappa_{K+8}}{2}, \kappa_{K+8} \right),$$

and  $\mathbf{w}$  is the  $3(K+7) \times 1$  vector given by

$$\mathbf{w} = \left\{ \frac{1}{6}(\Delta\boldsymbol{\kappa})_1, \frac{4}{6}(\Delta\boldsymbol{\kappa})_1, \frac{1}{6}(\Delta\boldsymbol{\kappa})_1, \frac{1}{6}(\Delta\boldsymbol{\kappa})_2, \frac{4}{6}(\Delta\boldsymbol{\kappa})_2, \right. \\ \left. \frac{1}{6}(\Delta\boldsymbol{\kappa})_2, \dots, \frac{1}{6}(\Delta\boldsymbol{\kappa})_{K+7}, \frac{4}{6}(\Delta\boldsymbol{\kappa})_{K+7}, \frac{1}{6}(\Delta\boldsymbol{\kappa})_{K+7} \right\},$$

where  $(\Delta\boldsymbol{\kappa})_k \equiv \kappa_{K+1} - \kappa_K$ ,  $1 \leq k \leq K+7$ . A common default choice for the number of knots is  $K = \min(n_U/4, 35)$ , where  $n_U$  is the number of unique  $x_i$ s, and the distribution of knots can either be quantile-based or equally spaced (e.g. Ruppert *et al.*, 2003). In the next section we show how the O'Sullivan penalised splines can be expressed in the mixed model and Bayesian hierarchical model framework.

### 1.3.1 Mixed model representation

Ruppert *et al.* (2003) presents the general framework of semiparametric regression using the inferential equivalence between penalised likelihood models and mixed models. In this subsection we explain why one can view penalised splines as mixed models. Consider the regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_u^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sigma_\varepsilon^2 \mathbf{I} \end{bmatrix} \right), \quad (1.4)$$

for the general design matrices  $\mathbf{X}$  and  $\mathbf{Z}$ . Let  $\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$ , the least squares estimator to (1.4) can be expressed as a best linear unbiased predictor given as

$$\hat{\boldsymbol{\nu}} = \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \left( \mathbf{C}^\top \mathbf{C} + \lambda \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \right)^{-1} \mathbf{C}^\top \mathbf{y}, \quad (1.5)$$

where the smoothing parameter is  $\lambda = \sigma_u^2 / \sigma_\varepsilon^2$ . The equivalence of (1.3) can be achieved if a  $(K+4) \times (K+4)$  linear transformation matrix  $\mathbf{L}$  can be found such that

$$\mathbf{C} = \mathbf{B}\mathbf{L} \quad \text{and} \quad \mathbf{L}^\top \boldsymbol{\Omega} \mathbf{L} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

According to Wand and Ormerod (2008), the usual method for obtaining  $\mathbf{L}$  and  $\boldsymbol{\Omega}$  is through spectral decomposition, and their resultant forms are respectively

$$\mathbf{L} = [\mathbf{U}_X \mid \mathbf{U}_Z \text{diag}(\mathbf{d}_Z^{-1/2})] \quad \text{and} \quad \boldsymbol{\Omega} = \mathbf{U} \text{diag}(\mathbf{d}) \mathbf{U}^\top,$$

where  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$  and  $\mathbf{d}$  is a  $(K+4) \times 1$  vector with exactly two zero entries and  $K+2$  positive entries. The term  $\mathbf{d}_Z$  is the  $(K+2) \times 1$  is a subvector of  $\mathbf{d}$  containing these positive entries and  $\mathbf{U}_Z$  is the  $(K+4) \times (K+2)$  submatrix of  $\mathbf{U}$  with columns corresponding to the positive entries of  $\mathbf{d}$ . Following this, model (1.4) with the O'Sullivan penalised splines can be fitted with the following design matrices:

$$\mathbf{X} = \mathbf{B}\mathbf{U}_X \quad \text{and} \quad \mathbf{Z} = \mathbf{B}\mathbf{U}_Z \text{diag}(\mathbf{d}_Z^{-1/2}).$$

## 1.4 Graphical models

Graphical models, originated from the field of machine learning and pattern recognition, are a branch of mathematics that combines ideas from graph theory and probability (e.g. Bishop, 2006; Jordan, 2004). Graphical models, also known as the *probabilistic graphical models*, are diagrammatic representations of probability distributions. A graph comprises



a set of *nodes* connected by *edges*. Adapting the conventions of Bishop (2006), nodes are shown as circles and edges are shown as line segments. Each node represents a random variable (or a set of random variables), and the edges express probabilistic relationships between these variables. There are two major types of graphical models: *directed acyclic graphs* (DAGs), also known as the *Bayesian networks*, and *undirected graphs*, also known as the *Markov random fields*. DAGs are distinguished from the undirected graphs by their directed edges, each of which has a single-headed arrow (double-headed arrows are not allowed). In both cases the graph expresses the way in which the joint density function over all of the random variables can be decomposed into a product of factors each depending on a subset of nodes, but the underlying relationship between the nodes and the factorisation is different for the two types of graphs. Of the two types, DAGs are more relevant to the theme of this thesis, and therefore we restrict our attention to DAGs.

Figure 1.1 is an elementary example of a DAG. In the graphical models literature, the variables  $x_1, x_2, x_3, x_5, x_6$  are random (“*hidden*”) nodes denoted by  $\mathcal{H}$ , corresponding to each of the open circles. The variable  $x_4$  is an observed (“*evidence*”) node denoted by  $\mathcal{E}$ , corresponding to the shaded circle. A DAG can be viewed as a “family tree”, where each directed edge signifies a “parent-child” relationship between the corresponding nodes. For example, the node  $x_2$  meets by the arrow-head is a *child* of the node  $x_4$  (*parent*). The nodes  $x_2$  and  $x_3$  have a common child  $x_4$  are called *co-parents* of each other. The *ancestors* of a node, say  $x_4$ , is the set consisting of the node’s parents, parents of the node’s parents, and so on, i.e.  $\{x_1, x_2, x_3, x_4\}$ . The *Markov blanket* of a node, say  $x_4$ , is the set consisting of the node’s parents, co-parents and children, i.e.  $\{x_2, x_3, x_5, x_6\}$ , as indicated by the blue dashed line polygon in Figure 1.1. The Markov blanket of a node separates that node from the remainder of the graph probabilistically and provides an important *locality property of probabilistic graphs* as we shall elaborate later on.

The DAG structure can easily be used to convey conditional independence properties of a model using the notions of *smallest ancestral subgraphs* and *moral graphs*. Let  $S = \{x_1, x_2, x_3, x_4\}$  be a subset of nodes, then an ancestral subgraph is a subgraph of a DAG, which, for any node in  $S$ , each of the node’s ancestors are also in the subgraph. The smallest ancestral subgraph containing  $S$  is the ancestral subgraph that has the fewest number of nodes (see Figure 1.2 left). The moral graph of a DAG is obtained by linking any pair of unconnected parents of a common child node with an undirected edge and changing all directed edges to undirected edges (referred to as *moralisation*, see Figure 1.2 right). Using the moral (undirected) graph in Figure 1.2, and suppose  $A = \{x_1\}$ ,  $B = \{x_4\}$  and  $C = \{x_2, x_3\}$ , then the conditional independence structure for these three sets of nodes is

$$A \perp\!\!\!\perp B \mid C \quad \text{i.e.} \quad x_1 \perp\!\!\!\perp x_4 \mid \{x_2, x_3\},$$

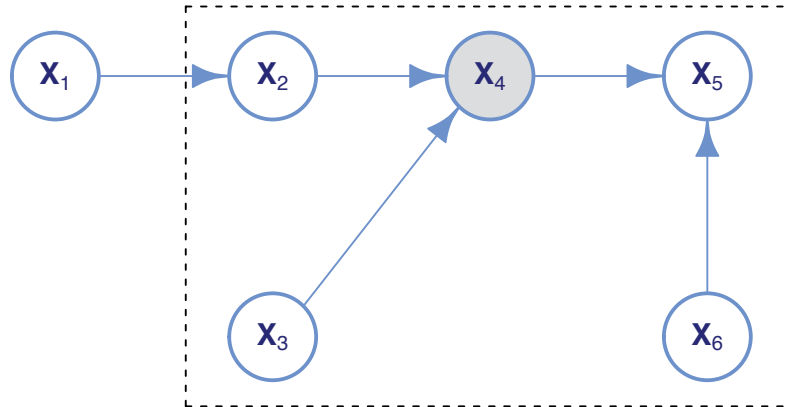


Figure 1.1: Directed acyclic graph involving six random variables:  $x_1, \dots, x_6$ . The open circles represent the random (“hidden”) nodes. The shaded circle represents the observed (“evidence”) node. The grey dashed line polygon represents the Markov blanket of the node  $x_4$ .

given that all the paths connecting nodes in  $A$  and nodes in  $B$  would pass through one or more nodes in  $C$ . The set of nodes in  $C$  blocks the paths from  $A$  to  $B$  and therefore the conditional independence holds. This conditional independence can also be derived using the *d-separation theorem* for a DAG.

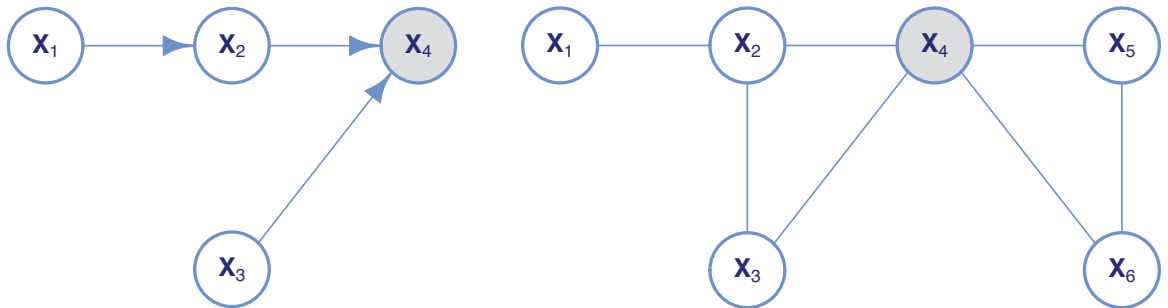


Figure 1.2: Left: The smallest ancestral subgraph corresponding to Figure 1.1. Right: The moral graph corresponding to Figure 1.1.

Through exploiting the graph-theoretic representation, the joint density function over all of the random variables can be factorised into a product of conditional density functions,

$$\begin{aligned}
 p(x_1, x_2, x_3, x_4, x_5, x_6) &= \prod_{i=1}^6 p(x_i \mid \text{parents of } x_i) \\
 &= p(x_1) p(x_2 \mid x_1) p(x_3) p(x_4 \mid x_2, x_3) p(x_5 \mid x_4, x_6) p(x_6),
 \end{aligned}$$

where parenthood is determined by the directed edges of the DAG. As demonstrated, conditional independence properties are useful in simplifying the structure of probabilis-

tic graphical models and thus more efficient algorithms can be devised for performing statistical inference. In a Bayesian setting, the model parameters are treated as random variables and represented as nodes in the graph. Computing the posterior distribution of those parameters is simply yet another inference problem!

In the next section, we describe probability and statistics from a Bayesian viewpoint, which provides a principled formalism through which all sources of uncertainty can be addressed consistently. Graphical models offers several useful properties in such a viewpoint including: (i) they provide a simple way to visualise the hierarchical structure of a probabilistic model; (ii) insights into conditional properties of the model can be obtained via inspection of the graph; and (iii) complex computations for performing inference can be expressed in terms of graphical manipulations and therefore the underlying algebraic expressions are carried along implicitly (Bishop, 2008).

## 1.5 Bayesian inference

Bayesian inference involves finding the joint posterior distribution of parameters of interest given the observed data. In many instances, exact inference is infeasible due to the posterior distributions being intractable, requiring approximate inference methods for use in practice. In recent years, Bayesian inference engines have emerged for approximate inference for general classes of mixed models. Examples include BUGS (Ligges *et al.*, 2009) and Stan (Stan Development Team, 2015), which are based on MCMC methods. In Subsection 1.5.2, we describe MCMC being one of the current standard methods for estimating Bayesian statistical models and has proven useful in a wide range of problems. However, for large models with complex posteriors, MCMC can be computationally intensive and suffers from poor mixing that leads to slow convergence. We therefore opt for an alternative approach, namely *variational approximations*, as we shall elaborate on in Subsection 1.5.3.

### 1.5.1 A Bayesian viewpoint of probability and statistics

Consider a generic Bayesian model with parameter vector  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_n\} \in \Theta$  and observed data vector  $\mathbf{y}$ . The goal of Bayesian inference is to infer the conditional density function of posterior  $\boldsymbol{\theta}$  given  $\mathbf{y}$  (known as the *posterior density function*),

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}, \boldsymbol{\theta})}{p(\mathbf{y})} = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}.$$

The numerator  $p(\mathbf{y}, \boldsymbol{\theta})$  is known as the *joint density function* of  $\mathbf{y}$  and  $\boldsymbol{\theta}$ , which can be written as a product of the sampling density function  $p(\mathbf{y}|\boldsymbol{\theta})$  and the prior density function  $p(\boldsymbol{\theta})$ . The denominator  $p(\mathbf{y})$  is known as the *marginal likelihood*, where  $p(\mathbf{y}) =$

$\int_{\Theta} p(\boldsymbol{\theta}) p(\mathbf{y}|\boldsymbol{\theta}) d\boldsymbol{\theta}$  in the case of continuous  $\boldsymbol{\theta}$ , or replace  $\int_{\Theta}$  with  $\sum_{\Theta}$  in the case of discrete  $\boldsymbol{\theta}$ .

For many models of practical interest the posterior density functions are intractable; they cannot be used to directly calculate the marginal density functions of parameters or other quantities of interest. For a long time the only generally applicable method was to use MCMC sampling techniques. MCMC is a sampling-based simulation algorithm which circumvents the intractability problem by generating a Markov Chain of samples of  $\boldsymbol{\theta}^{(t)}$ ,  $0 \leq t \leq T$ , to approximate the target posterior density function  $p(\boldsymbol{\theta}|\mathbf{y})$ . A *Markov Chain* is a sequence of random variables  $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(T)}$  for which

$$p(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(t-1)}, \mathbf{y}) = p(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t-1)}, \mathbf{y}),$$

so that the current state depends only on the preceding state.

### 1.5.2 Markov chain Monte Carlo

Gibbs sampling is the simplest component-wise algorithm of the MCMC methods (Gelfand and Smith, 1990; Geman and Geman, 1984; Robert and Casella, 2004). The Markov chain of samples of  $\boldsymbol{\theta}^{(t)}$ ,  $0 \leq t \leq T$ , is obtained by directly sampling from the *full conditional density functions*  $p\left(\theta_i^{(t)}|\theta_{[i]}^{(t-1,t)}\right)$ , specifying the density of  $\theta_i$  conditional on other parameters  $\theta_{[i]}$ . Each subvector  $\theta_i^{(t)}$  is updated conditional on the latest value of the other parameters  $\theta_{[i]}$ , which are the iteration  $t$  values for the parameters already updated and the iteration  $t - 1$  values for the others. This is illustrated in Algorithm 1.

---

**Initialise:**  $\theta_1^{(0)}, \dots, \theta_n^{(0)}$

**Cycle:**

For  $t = 0, \dots, T$ :

$$\theta_1^{(t)} \sim p(\theta_1|\theta_2^{(t-1)}, \dots, \theta_n^{(t-1)})$$

$$\theta_2^{(t)} \sim p(\theta_2|\theta_1^{(t)}, \dots, \theta_n^{(t-1)})$$

$$\vdots$$

$$\theta_{n-1}^{(t)} \sim p(\theta_{n-1}|\theta_1^{(t)}, \dots, \theta_n^{(t-1)})$$

$$\theta_n^{(t)} \sim p(\theta_n|\theta_1^{(t)}, \dots, \theta_{n-1}^{(t)})$$

**until convergence reaches.**

---

Algorithm 1: The Gibbs sampler algorithm in generic form.

Unfortunately, this approach tends to be computationally costly for complex models

that leads to unacceptably slow convergence and does not scale well to real-world applications involving large datasets.

To address this problem, a research group from Columbia University has created **Stan**: new, high-performance open-source software for Bayesian inference on multilevel models (Stan Development Team, 2015). Rather than the conventional Gibbs sampler, **Stan** uses a variant of Hamiltonian Monte Carlo algorithm to speed up convergence of generalised multilevel linear models. Instead of the typical random-walk behaviours, Hamiltonian Monte Carlo is a MCMC method that simulates a physical system governed by Hamiltonian dynamics (a way that physicists use to describe an object’s motion in terms of its location and momentum at some time point) to propose future states in the Markov chain. This allows the Markov chain to explore the posterior distribution much more efficiently, resulting in faster convergence.

### 1.5.3 A brief introduction to variational approximations

Variational approximations are a fast, versatile and deterministic class of approximations with origins in the statistical physics and computer science literature (Bishop, 2006; Jordan *et al.*, 1999; Titterton, 2004; Wainwright and Jordan, 2008). Since about 2005, variational methods have been increasingly explored in the statistical literature. For example, McGrory and Titterton (2007), Wand *et al.* (2012) and Ormerod and Wand (2010) present variational methodology for a diversified range of applications, from finite mixture models, complex models with elaborate distributions (such as asymmetric Laplace and skew normal) to spline and wavelet regression models. In addition, Wang and Titterton (2006) prove convergence of variational algorithms for normal mixture models and You *et al.* (2014) propose several information criteria that are useful for model selection. We provide here a brief overview of variational approximations and highlight the key concepts. Comprehensive summaries can be found in Bishop (2006) and Ormerod and Wand (2010).

#### 1.5.3.1 Mean field variational Bayes

Consider a generic Bayesian model with observed data vector  $\mathbf{y}$  and parameter vector  $\boldsymbol{\theta}$  that is continuous over the parameter space  $\Theta$ . The essence of variational inference is to approximate the posterior density function  $p(\boldsymbol{\theta}|\mathbf{y})$  with a so-called *approximating density function*  $q(\boldsymbol{\theta})$ . To make this approximation as close as possible, we search over  $q \in \mathcal{Q}$ , for some set  $\mathcal{Q}$  of density functions, to find a particular density function with the minimum Kullback-Liebler (KL) distance/divergence with the actual posterior

$$q^*(\boldsymbol{\theta}) = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \{q(\boldsymbol{\theta}) \| p(\boldsymbol{\theta}|\mathbf{y})\}, \quad (1.6)$$

where

$$\text{KL} \{q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}|\mathbf{y})\} = \int_{\Theta} q(\boldsymbol{\theta}) \log \left\{ \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y})} \right\} d\boldsymbol{\theta}.$$

The approximating distribution can be chosen to have a simple analytical form. For example, it could be a Gaussian distribution, whose mean and covariance are optimised to minimise the Kullback-Liebler distance with respect to the actual posterior. A more flexible framework is to impose some form of factorisation or product density restriction on  $q(\boldsymbol{\theta})$ , without any restriction on the functional form of the factors (Bishop, 2006):

$$\mathcal{Q} = \left\{ q(\boldsymbol{\theta}) : q(\boldsymbol{\theta}) = \prod_{i=1}^M q_i(\boldsymbol{\theta}_i) \text{ for some partition } \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\} \text{ of } \boldsymbol{\theta} \right\}, \quad (1.7)$$

known as the *mean field restriction*. The resultant  $q^*(\boldsymbol{\theta})$  is known as the *mean field variational Bayes approximation* (MFVB) to the actual posterior, and from now on we refer to it as *optimal  $q$ -density*. Essentially, it assumes posterior independence among parameters that may not be present in the actual posterior. Depending on the chosen partition of  $\boldsymbol{\theta}$ , this independence assumption may be rather unrealistic in settings where there exists high posterior correlations among parameters, thus leading to poor approximations. The particular parametric families that constitute each of the approximating  $q$ -density factors  $q_i(\boldsymbol{\theta}_i)$  are derived through the variational methodology. Indeed, minimising the Kullback-Liebler divergence in (1.6) is equivalent to maximising the lower bound  $\underline{p}(\mathbf{y}; q)$  since the marginal log-likelihood can be expressed as (Ormerod and Wand, 2010)

$$\begin{aligned} \log p(\mathbf{y}) &= \log p(\mathbf{y}) \int_{\Theta} q(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\Theta} q(\boldsymbol{\theta}) \log p(\mathbf{y}) d\boldsymbol{\theta} \\ &= \int_{\Theta} q(\boldsymbol{\theta}) \log \left\{ \frac{p(\mathbf{y}, \boldsymbol{\theta})/q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y})/q(\boldsymbol{\theta})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} q(\boldsymbol{\theta}) \log \left\{ \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} \right\} d\boldsymbol{\theta} + \int_{\Theta} q(\boldsymbol{\theta}) \log \left\{ \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y})} \right\} d\boldsymbol{\theta} \\ &= \log \underline{p}(\mathbf{y}; q) + \text{KL}\{q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}|\mathbf{y})\} \end{aligned}$$

with  $\text{KL} \{q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}|\mathbf{y})\} \geq 0$  for all densities  $q$ .

The optimal  $q$ -density functions under product restriction (1.7) can be obtained via an *iterative coordinate ascent algorithm* that is analogous to the Expectation-Maximisation algorithm (Bishop, 2006; Ormerod and Wand, 2010). Define

$$E_{q(-\boldsymbol{\theta}_i)} \{\log p(\mathbf{y}, \boldsymbol{\theta})\} = \int \log p(\mathbf{y}, \boldsymbol{\theta}) \prod_{j \neq i} q(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j$$

to be the log posterior averaged over the current estimates of the approximating density functions for all but the  $i$ th parameter vector. The term  $E_{q(-\boldsymbol{\theta}_i)}$  denotes expectation with respect to the  $q$ -densities of all parameters except  $\boldsymbol{\theta}_i$ . Variational approximations is now reduced to solving an optimisation problem in the form of (1.6). We proceed by initialising each of the  $q$ -density factors  $q_i(\boldsymbol{\theta}_i)$  and updating each factor successively using the current estimates of the other factors. At the end of each iteration, an updated value of the lower bound is computed,

$$\log \underline{p}(\mathbf{y}; q) = E_q\{\log p(\mathbf{y}, \boldsymbol{\theta}) - \log q(\boldsymbol{\theta})\}, \quad (1.8)$$

and the algorithm is iterated until convergence of the lower bound to its maximum, see Algorithm 2. Convergence of such simple iterative updates is guaranteed in most cases since the algorithm is fundamentally a generalisation of Expectation-Maximisation (Chappell *et al.*, 2009). The second term on the right hand side of (1.8) corresponds to the *entropy* of a density function. If  $x$  is a random vector with a density function  $p$  then the corresponding entropy is given by

$$\text{Entropy}(p) \equiv E_p\{-\log p(x)\}.$$

For many common distribution families, the entropy can be expressed algebraically in terms of the distribution's parameters. For example, if  $p(x; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the multivariate normal density function of dimension  $d$  with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  then

$$\text{Entropy}(p; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2} d \{1 + \log(2\pi)\} + \frac{1}{2} \log|\boldsymbol{\Sigma}|.$$

Another common entropy expression that arises throughout this thesis is the inverse-Gamma family of density functions. The entropy corresponds to  $p(x; \kappa, \lambda)$  with parameters  $\kappa, \lambda > 0$  is

$$\text{Entropy}(p; \kappa, \lambda) = \log(\lambda) + \kappa + \log\{\Gamma(\kappa)\} - (\kappa + 1) \text{digamma}(\kappa),$$

where  $\text{digamma}(\kappa) \equiv (d/dx) \log \Gamma(\kappa)$  is the digamma function.

While the distributional forms of the approximating density functions  $q_i(\boldsymbol{\theta}_i)$  are unspecified, the structure of the statistical model itself lends to a solution that lies in a particular parametric family for each of the  $q_i(\boldsymbol{\theta}_i)$ . For example, when all the parameters in a model are conditionally conjugate, the optimal  $q$ -density functions in Algorithm 2 are available in closed form. If  $q_i(\boldsymbol{\theta}_i)$  can not be recognised as a standard distribution, then numerical integration methods are required to estimate the marginal likelihood, which is computationally more demanding.

**Initialise:**  $q_1(\boldsymbol{\theta}_1), \dots, q_M(\boldsymbol{\theta}_M)$

**Cycle:**

$$q_1(\boldsymbol{\theta}_1) \leftarrow \frac{\exp[E_{q(-\boldsymbol{\theta}_1)}\{\log p(\mathbf{y}, \boldsymbol{\theta})\}]}{\int \exp[E_{q(-\boldsymbol{\theta}_1)}\{\log p(\mathbf{y}, \boldsymbol{\theta})\}]d\boldsymbol{\theta}_1}$$

$$\vdots$$

$$q_M(\boldsymbol{\theta}_M) \leftarrow \frac{\exp[E_{q(-\boldsymbol{\theta}_M)}\{\log p(\mathbf{y}, \boldsymbol{\theta})\}]}{\int \exp[E_{q(-\boldsymbol{\theta}_M)}\{\log p(\mathbf{y}, \boldsymbol{\theta})\}]d\boldsymbol{\theta}_M$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

---

Algorithm 2: Iterative scheme for obtaining the optimal  $q$ -density functions under product restriction (1.7).

Upon convergence, Algorithm 2 shows that the optimal occurs when

$$q_i^*(\boldsymbol{\theta}_i) \propto \exp [E_{(-\boldsymbol{\theta}_i)}\{\ln p(\boldsymbol{\theta}, \mathbf{y})\}], \quad 1 \leq i \leq M.$$

This means that the log of the optimal solution for each factor is simply obtained by considering the log of the joint density function over all of the random and observed variables and taking the expectation with respect to all of the other factors  $q_j(\boldsymbol{\theta}_j)$ ,  $j \neq i$  (Ormerod and Wand, 2010). An alternative expression for  $q_i^*(\boldsymbol{\theta}_i)$  is

$$q_i^*(\boldsymbol{\theta}_i) \propto \exp [E_{q(-\boldsymbol{\theta}_i)}\{\ln p(\boldsymbol{\theta}_i|\text{rest})\}], \quad 1 \leq i \leq M, \quad (1.9)$$

where “rest” denotes all of the random variables excluding  $\boldsymbol{\theta}_i$ , i.e.  $\{\mathbf{y}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{i-1}, \boldsymbol{\theta}_{i+1}, \dots, \boldsymbol{\theta}_M\}$ , and the distribution  $p(\boldsymbol{\theta}_i|\text{rest})$  is known as the *full conditional density function*. The formulation of (1.9) greatly benefits from a DAG-related concept known as the *Markov blanket theory* - Theorem 1.5.1. The Markov blanket of a node on a DAG is the set of children, parents and coparents of that node. Two nodes are coparents if they have at least one child node in common.

**Theorem 1.5.1.** For each node on a probabilistic DAG, the conditional distribution of the node given the rest of the nodes is the same as the conditional distribution of the node given its Markov blanket (Dechter and Pearl, 1988).

For our generic Bayesian example, this implies that

$$p(\boldsymbol{\theta}_i|\text{rest}) = p(\boldsymbol{\theta}_i|\text{Markov blanket of } \boldsymbol{\theta}_i).$$



It immediately follows that

$$q_i^*(\boldsymbol{\theta}_i) \propto \exp [E_{q(-\boldsymbol{\theta}_i)}\{\log p(\boldsymbol{\theta}_i|\text{Markov blanket of } \boldsymbol{\theta}_i)\}], \quad 1 \leq i \leq M. \quad (1.10)$$

This is known as the *locality property of DAGs*. For large DAGs, this property yields considerable algebraic benefits. In particular, it shows that the  $q_i^*(\boldsymbol{\theta}_i)$  require only local calculations on the model's DAG.

### 1.5.3.2 Variational message passing

*Variational message passing* or VMP (e.g. Winn and Bishop, 2005) is a message passing algorithmic implementation of the mean field approximation, limited to conjugate-exponential models. The VMP works with *messages* passed between neighbouring nodes (i.e. parent and child) on a *factor graph* of the model. Figure 1.3 shows an example of a factor graph corresponding to the DAG in Figure 1.1. Examples of the messages are:

$$m_{p(x_1) \rightarrow x_1}(x_1), \quad m_{x_1 \rightarrow p(x_1)}(x_1), \quad m_{x_1 \rightarrow p(x_2|x_1)}(x_1), \quad m_{p(x_2|x_1) \rightarrow x_1}(x_2).$$

The VMP messages are equivalent to the q-density factors  $q_i(\boldsymbol{\theta}_i)$  in MFVB, and their conditional density functions are members of the exponential family including Gaussian, Poisson, Gamma etc. Thus, calculating  $q_i(\boldsymbol{\theta}_i)$  simply involves summing sufficient statistics, and similarly, the variational optimisation can be decomposed into a set of local computations that depend only on messages from neighbouring nodes in the DAG.

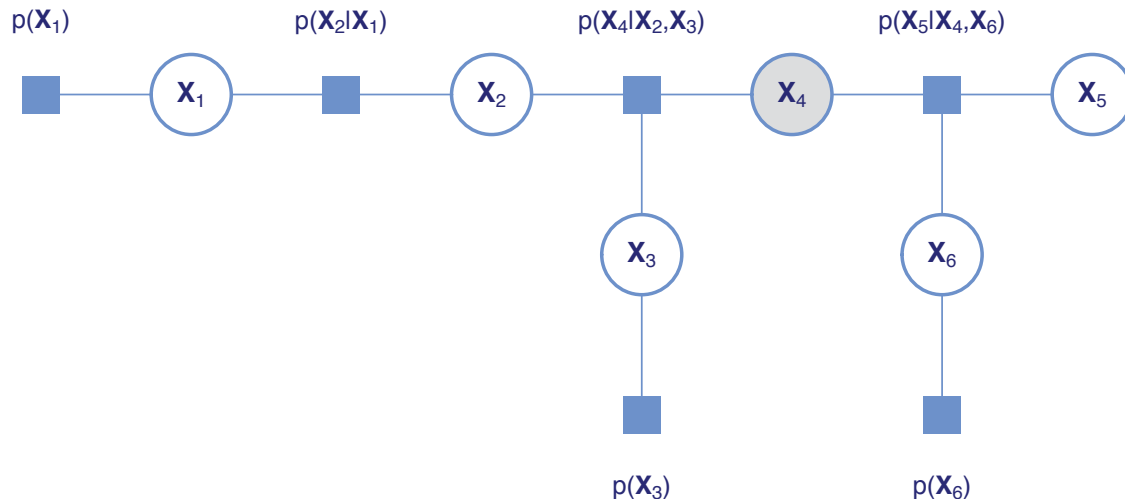


Figure 1.3: Factor graph corresponding to the DAG in Figure 1.1.

An in-depth explanation as well as a discussion on the factor graphs and VMP is deferred to Chapter 6.

### 1.5.3.3 Non-conjugate variational message passing

Knowles and Minka (2011) introduce *non-conjugate variational message passing* (NCVMP) which widens the scope of tractable models for VMP and MFVB in general. Specifically, it extends variational inference to non-conjugate exponential models. The modularity of NCVMP allows modifications only to those model parameters that involve previously intractable solutions, whilst keeping VMP or MFVB solutions of the other model parameters the same. Let us reconsider the generic Bayesian model with parameter vector  $\boldsymbol{\theta}$  and observed data vector  $\mathbf{y}$ , but this time include an additional parameter vector  $\boldsymbol{\phi}$ . From (1.7), the MFVB approximates the joint posterior density function  $p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{y})$  by the factorised form

$$q_{\boldsymbol{\theta}_1}(\boldsymbol{\theta}_1), \dots, q_{\boldsymbol{\theta}_M}(\boldsymbol{\theta}_M) q_{\boldsymbol{\phi}}(\boldsymbol{\phi}), \quad (1.11)$$

and suppose that  $q_{\boldsymbol{\theta}_1}(\boldsymbol{\theta}_1), \dots, q_{\boldsymbol{\theta}_M}(\boldsymbol{\theta}_M)$  are in the exponential family and  $q_{\boldsymbol{\phi}}(\boldsymbol{\phi})$  is not. The optimal  $q$ -density functions satisfy

$$\begin{aligned} q_{\boldsymbol{\theta}_i}^*(\boldsymbol{\theta}_i) &\propto \exp [E_{q(-\boldsymbol{\theta}_i)} \{\log p(\boldsymbol{\theta}_i | \text{rest})\}], \quad 1 \leq i \leq M, \\ \text{and } q_{\boldsymbol{\phi}}^*(\boldsymbol{\phi}) &\propto \exp [E_{q(-\boldsymbol{\phi})} \{\log p(\boldsymbol{\phi} | \text{rest})\}]. \end{aligned}$$

The term  $E_{q(-\boldsymbol{\phi})} \{\log p(\boldsymbol{\phi} | \text{rest})\}$  is intractable due to difficulties arising from non-conjugacy. Non-conjugate variational message passing offers a remedy by postulating  $q_{\boldsymbol{\phi}}(\boldsymbol{\phi})$  to be an exponential family density function with the *natural parameter vector*  $\boldsymbol{\eta}$ , *natural statistic vector*  $\mathbf{T}(\boldsymbol{\phi})$  and the *log-partition function*  $A(\boldsymbol{\eta})$ :

$$q(\boldsymbol{\theta}; \boldsymbol{\eta}) = \exp \left\{ \mathbf{T}(\boldsymbol{\phi})^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) + \text{const} \right\}.$$

where “const” denotes additive constants with respect to the function argument. Hence, (1.11) becomes

$$q_{\boldsymbol{\theta}_1}(\boldsymbol{\theta}_1), \dots, q_{\boldsymbol{\theta}_M}(\boldsymbol{\theta}_M) q_{\boldsymbol{\phi}}(\boldsymbol{\phi}; \boldsymbol{\eta}).$$

From Theorem 1 of Knowles and Minka (2011), the optimal  $q$ -density functions can then be found using

$$\begin{aligned} q_{\boldsymbol{\theta}_i}^*(\boldsymbol{\theta}_i) &\propto \exp [E_{q(-\boldsymbol{\theta}_i)} \{\log p(\boldsymbol{\theta}_i | \text{rest})\}], \quad 1 \leq i \leq M, \\ \boldsymbol{\eta} &\leftarrow [\text{var} \{\mathbf{T}(\boldsymbol{\phi})\}]^{-1} \{D_{\boldsymbol{\eta}} E_{q(\boldsymbol{\theta}, \boldsymbol{\phi})} \{\log p(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\phi})\}\}, \end{aligned} \quad (1.12)$$

where  $D_{\mathbf{x}}f$  is defined in Section 1.10. This thesis only concerns a special case where  $q(\boldsymbol{\theta}; \boldsymbol{\eta})$  corresponds to a  $d$ -dimensional multivariate normal density function. The advantage of

exponential family density functions is that expectations of their logarithms are tractable to compute and their state can be summarised completely by the natural parameter vector. A recent example of utilising non-conjugate VMP for approximate Bayesian inference in generalised linear mixed models is given in Tan and Nott (2013).

## 1.6 Terminological and notational conventions

Lower-case Roman and Greek letters in boldface are used to denote random vectors with entries consisting of subscripts. For example,  $\mathbf{x}$  denotes a  $n \times 1$  vector containing  $x_1, \dots, x_n$ . Upper-case Roman and Greek letters in boldface are used to denote random matrices. For example,  $\mathbf{X}$  denotes a  $m \times n$  matrix containing  $n$  vectors of dimension  $m \times 1$ , i.e.  $\mathbf{x}_1, \dots, \mathbf{x}_m$ . Scalar functions applied to vectors are evaluated element-wise. For example,

$$\exp(a_1, a_2, a_3) \equiv (\exp(a_1), \exp(a_2), \exp(a_3)).$$

For two matrices  $\mathbf{A}$  and  $\mathbf{B}$  of the same dimensions, the element-wise product denoted  $\mathbf{A} \odot \mathbf{B}$  is a matrix with elements  $(\mathbf{A} \odot \mathbf{B})_{ij} = (\mathbf{A})_{ij} (\mathbf{B})_{ij}$ . For two matrices  $\mathbf{A}$  and  $\mathbf{B}$  of any different dimension  $m \times n$  and  $p \times q$  respectively, the Kronecker product denoted  $\mathbf{A} \otimes \mathbf{B}$  is a matrix with dimension  $mn \times pq$  with elements  $(\mathbf{A} \otimes \mathbf{B})_{ij} = (\mathbf{A})_{ij} \mathbf{B}$ . If  $\mathbf{A}$  and  $\mathbf{B}$  have the same number of rows, then the notation  $(\mathbf{A}|\mathbf{B})$  denotes concatenation (by columns) of matrices  $\mathbf{A}$  and  $\mathbf{B}$ . We use  $\mathbf{1}_d$  to denote the  $d \times 1$  vector with all entries equal to 1 and  $\mathbf{I}_d$  to denote the  $d \times d$  identity matrix. The norm of a vector denoted  $\|\mathbf{v}\|$  is defined to be  $\sqrt{\mathbf{v}^\top \mathbf{v}}$ . For a  $n \times n$  symmetric matrix  $\mathbf{M}$ , the trace and determinant are denoted by  $\text{tr}(\mathbf{M})$  and  $|\mathbf{M}|$  respectively and adopt the usual definitions.

For a  $d \times 1$  vector  $\mathbf{a}$ , we use  $\text{diag}(\mathbf{a})$  to denote the  $d \times d$  diagonal matrix containing entries of  $\mathbf{a}$  along the main diagonal. For a  $d \times d$  square matrix  $\mathbf{A}$ , we use  $\text{diagonal}(\mathbf{A})$  to denote the  $d \times 1$  vector containing the main diagonal entries of  $\mathbf{A}$ . For square matrices  $\mathbf{A}_1, \dots, \mathbf{A}_r$ , we use  $\text{blockdiag}(\mathbf{A}_1, \dots, \mathbf{A}_r)$  to denote the block diagonal matrix with the  $i^{\text{th}}$  block equal to  $\mathbf{A}_i$ . The transpose of  $\mathbf{A}$  is denoted by  $\mathbf{A}^\top$ . We let  $\mathbf{A}^{-1}$  denote the inverse of a matrix provided  $\mathbf{A}$  is an invertible matrix.

For a  $d \times d$  matrix  $\mathbf{A}$ , we use  $\text{vec}(\mathbf{A})$  to denote the  $d^2 \times 1$  vector obtained by concatenating the columns of  $\mathbf{A}$  underneath each other from left to right. The term  $\text{vech}(\mathbf{A})$  denotes the  $\frac{1}{2}d(d+1) \times 1$  vector obtained from  $\text{vec}(\mathbf{A})$  by eliminating the above-diagonal entries of  $\mathbf{A}$ . For example,

$$\text{vec} \left( \begin{bmatrix} 3 & 8 \\ 4 & 9 \end{bmatrix} \right) = \begin{bmatrix} 3 \\ 4 \\ 8 \\ 9 \end{bmatrix} \quad \text{and} \quad \text{vech} \left( \begin{bmatrix} 3 & 8 \\ 4 & 9 \end{bmatrix} \right) = \begin{bmatrix} 3 \\ 4 \\ 9 \end{bmatrix}.$$

Let  $u$  and  $v$  be two continuous or discrete random scalar variables, we denote  $p(u)$  as the density function (probability mass function for discrete and probability density function for continuous). This density function notation extends, analogously, to joint and conditional density functions. For example,  $p(u, v)$  denotes the joint density function of  $u$  and  $v$ , and  $p(u|v)$  denotes the conditional density function of  $u$  given by  $v$ . The expectation and covariance matrix of  $u$  are denoted by  $E(u)$  and  $\text{Var}(u)$  respectively. Analogous terminology applies to random vectors and matrices.

The  $q$  denotes the density functions that arise from the MFVB approximation. For a generic random scalar variable  $\theta$  and density function  $q$  we define

$$\mu_{q(\theta)} \equiv E_q(\theta) \quad \text{and} \quad \sigma_{q(\theta)} \equiv \text{Var}_q(\theta).$$

For a generic random vector  $\boldsymbol{\theta}$  and density function we define

$$\boldsymbol{\mu}_{q(\boldsymbol{\theta})} \equiv E_q(\boldsymbol{\theta}) \quad \text{and} \quad \boldsymbol{\Sigma}_{q(\boldsymbol{\theta})} \equiv \text{Cov}_q(\boldsymbol{\theta}).$$

## 1.7 Common distributions

If  $y_i$  has a distribution  $D_i$  for each  $1 \leq i \leq n$  and the  $y_i$  are independent, then we write  $y_i \stackrel{\text{ind.}}{\sim} D_i$ . Table 1.1 lists all of the distributions used in this thesis. In particular, the parametrisation of the corresponding density functions is provided.

## 1.8 Distributional results

**Result 1.8.1.** Let  $x$  and  $a$  be random variables such that

$$x|a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/a) \quad \text{and} \quad a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/A^2).$$

Then  $\sqrt{x} \sim \text{Half-Cauchy}(A)$ .

The above hierarchical representation enables MCMC and variational methods easier to carry out because of the conditional conjugacy properties of the Inverse-Gamma distribution. Extending Result 1.8.1 leads to the following:

$$\begin{aligned} \boldsymbol{\Sigma}|a_1, \dots, a_p &\sim \text{Inverse-Wishart}(\nu + p - 1, 2\nu \text{diag}(1/a_1, \dots, 1/a_p)) \\ \text{and } a_k &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\frac{1}{2}, 1/A_k^2), \quad k = 1, \dots, p, \end{aligned}$$

where  $\boldsymbol{\Sigma}$  is a  $p \times p$  random matrix and  $\mathbf{a}$  is a  $p \times 1$  random vector. According to Huang and Wand (2013) simulation results, the scale-mixture  $\boldsymbol{\Sigma}$  offers an overall increase in practical performance over the classical inverse-Wishart prior for sparse covariance matrix

1.8. DISTRIBUTIONAL RESULTS

Distribution	Probability density or mass function in $x$	Abbreviation
Bernoulli	$p^x (1-p)^{1-x}; \quad x = 0, 1; p \in (0, 1)$	Bernoulli( $p$ )
Poisson	$\lambda^x e^{-\lambda}/x!; \quad x = 0, 1, \dots; \lambda > 0$	Poisson( $\lambda$ )
Uniform	$1/(b-a); \quad a < x < b$	Uniform( $a, b$ )
Normal	$\frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}; \quad x, \mu \in \mathbb{R}; \sigma^2 > 0$	$N(\mu, \sigma^2)$
Student- $t$	$\frac{\Gamma\{(\nu+1)/2\}}{\sqrt{\pi\nu\sigma^2} \Gamma(\frac{\nu}{2}) \left\{1 + \frac{(x-\mu)^2}{\nu\sigma^2}\right\}^{(\nu+1)/2}}; \quad \sigma, \nu > 0$	$t(\mu, \sigma^2, \nu)$
Gamma	$\frac{B^A x^{A-1} e^{-Bx}}{\Gamma(A)}; \quad x > 0; A, B > 0$	Gamma( $A, B$ )
Inverse-Gamma	$\frac{B^A x^{-A-1} e^{-B/x}}{\Gamma(A)}; \quad x > 0; A, B > 0$	Inverse-Gamma( $A, B$ )
Half-Cauchy	$\frac{2\sigma}{\pi(x^2 + \sigma^2)}; \quad x > 0; \sigma > 0$	Half-Cauchy( $\sigma$ )
Multivariate Normal	$ 2\pi\mathbf{\Sigma} ^{-1/2} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$ $\mathbf{x} \in \mathbb{R}^d; \mathbf{\Sigma}$ is symmetric and positive definite	$N(\boldsymbol{\mu}, \mathbf{\Sigma})$
Inverse-Wishart	$(C_{d,A}^{-1}  \mathbf{B} ^{A/2}  \mathbf{X} ^{(A+d+1)/2} \exp\{-\frac{1}{2}\text{tr}(\mathbf{B}\mathbf{X})\})$ where $(C_{d,A} \equiv 2^{Ad/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma(\frac{A+1-j}{2}))$ . $A > 0; \mathbf{B}$ is symmetric and positive definite	Wishart( $A, \mathbf{B}$ )

Table 1.1: Distributions used in this thesis and their corresponding probability density or mass functions.

estimation. In addition, a compelling feature of the scale-mixture  $\boldsymbol{\Sigma}$  is that, its marginal distribution is invariant even when the matrix itself is expanded or collapsed over any set of variables in a self-consistent manner.

**Result 1.8.2.** Let  $x$  and  $b$  be random variables such that

$$x|b \sim N(\mu, b\sigma^2) \quad \text{and} \quad b \sim \text{Inverse-Gamma}(\frac{\nu}{2}, \nu/2).$$

Then  $\sqrt{x} \sim t(\mu, \sigma, \nu)$ .

**Result 1.8.3.** If  $\nu \sim \text{Inverse-Gamma}(A, B)$  then

$$E(1/\nu) = A/B \quad \text{and} \quad E\{\log(\nu)\} = \log(B) - \text{digamma}(A).$$

## 1.9 Vector differential calculus

**Definition 1.9.1.** The derivative factor of a scalar-valued function  $f$  with respect to  $x \in \mathbb{R}^p$  is denoted by  $D_x f$ , with the  $i$ th entry equals to

$$\frac{\partial f(\mathbf{x})_i}{\partial x_i}.$$

**Definition 1.9.2.** The Hessian matrix of a scalar-valued function  $f$  with respect to  $x \in \mathbb{R}^p$  is denoted by  $H_x f$ , with the  $(i, j)$ th entry equals to

$$\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}.$$

**Definition 1.9.3.** Let  $\mathbf{U}$  and  $\mathbf{V}$  are matrix functions, then

$$\begin{aligned} \text{vec}(d\mathbf{U}) &= d\text{vec}(\mathbf{U}) \\ d|\mathbf{U}| &= |\mathbf{U}| \text{tr}(\mathbf{U}^{-1} d\mathbf{U}) \\ d\mathbf{U}^{-1} &= -\mathbf{U}^{-1} (d\mathbf{U}) \mathbf{U}^{-1} \\ d(\text{tr}(\mathbf{U})) &= \text{tr}(d\mathbf{U}) \\ d(\mathbf{U} \odot \mathbf{V}) &= (d\mathbf{U}) \odot \mathbf{V} + \mathbf{U} \odot (d\mathbf{V}) \\ d(\log|\mathbf{U}|) &= \text{tr}(\mathbf{U}^{-1}) d\mathbf{U} \\ d(\mathbf{U}\mathbf{V}) &= (d\mathbf{U})\mathbf{V} + \mathbf{U}(d\mathbf{V}) \end{aligned}$$

## 1.10 Special functions

**Definition 1.10.1.** The integral  $\mathcal{F}$  is a non-analytic integral family, defined as

$$\mathcal{F}(p, q, r, s, t) \equiv \int_s^t x^p \exp \left[ q \{ \log(x/2) - \log \Gamma(x/2) \} - \frac{1}{2} r x \right] dx, \quad p \geq 0, q, r, s, t > 0.$$

**Definition 1.10.2.** The logit( $\cdot$ ) function is defined by

$$\text{logit}(p) = \log \left( \frac{p}{1-p} \right) = \log(p) - \log(1-p), \quad 0 \leq p \leq 1.$$

**Definition 1.10.3.** The digamma function is denoted by  $\psi(\cdot)$  and is the derivative of the logarithm of the gamma function

$$\psi(x) = \frac{d}{dx} \ln \Gamma(x) = \frac{\Gamma'(x)}{\Gamma(x)}, \quad x > 0.$$

**Definition 1.10.4.** The duplication matrix  $\mathbf{D}_d$  is an unique  $d^2 \times \frac{1}{2} d(d+1)$  matrix of

zeros and one that transforms any  $d \times d$  symmetric matrix  $\mathbf{A}$  from  $\text{vech}(\mathbf{A})$  into  $\text{vec}(\mathbf{A})$ :

$$\mathbf{D}_d \text{vech}(\mathbf{A}) = \text{vec}(\mathbf{A}) \quad \text{for } \mathbf{A} = \mathbf{A}^\top.$$

**Definition 1.10.5.** The Moore-Penrose Inverse of  $\mathbf{D}_d$ , defined as

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

is an unique  $\frac{1}{2}d(d+1) \times d^2$  matrix that transforms any  $d \times d$  symmetric matrix  $\mathbf{A}$  from  $\text{vec}(\mathbf{A})$  into  $\text{vech}(\mathbf{A})$ :

$$\mathbf{D}_d^+ \text{vec}(\mathbf{A}) = \text{vech}(\mathbf{A}).$$

**Definition 1.10.6.** The two identities involving  $\text{vec}$  and  $\text{vech}$  are:

$$\text{vec}^{-1}(\text{vec}(\mathbf{A})) = \mathbf{A} \quad \text{and} \quad \text{vec}(\text{vec}^{-1}(\mathbf{a})) = \mathbf{a},$$

where  $\mathbf{A}$  is a  $d \times d$  matrix,  $\mathbf{a}$  is a  $d^2 \times 1$  vector and  $\text{vec}^{-1}(\mathbf{a})$  is a  $d \times d$  matrix obtained from unstacking the entries of  $\mathbf{a}$  in a column-wise fashion.

## 1.11 Matrix results

**Result 1.11.1.** If  $\mathbf{A}$  is a square matrix, then  $\text{tr}(\mathbf{A}^\top) = \text{tr}(\mathbf{A})$ .

**Result 1.11.2.** If  $\mathbf{A}$  is a  $m \times n$  matrix and  $\mathbf{B}$  is a  $n \times m$  matrix, then  $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ .

**Result 1.11.3.** If  $\mathbf{A}$  and  $\mathbf{B}$  are  $m \times n$  matrices, then  $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$ .

**Result 1.11.4.** If  $\mathbf{a}$  and  $\mathbf{b}$  are vectors of the same length, then

$$\|\mathbf{a} - \mathbf{b}\|^2 = \|\mathbf{a}\|^2 - 2\mathbf{a}^\top \mathbf{b} + \|\mathbf{b}\|^2.$$

**Result 1.11.5.** Let  $\mathbf{A}$  be a symmetric invertible matrix and  $\mathbf{x}$  and  $\mathbf{b}$  be column vectors with the same number of rows as  $\mathbf{A}$ . Then

$$\frac{1}{2}\mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x} = -\frac{1}{2}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b})^\top \mathbf{A}(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}) + \frac{1}{2}\mathbf{b}^\top \mathbf{A}^{-1}\mathbf{b}.$$

**Result 1.11.6.** Let  $\mathbf{v}$  be a random vector. Then

$$E(\mathbf{v}\mathbf{v}^\top) = E(\mathbf{v})E(\mathbf{v})^\top + \text{Cov}(\mathbf{v}) \quad \text{and} \quad E(\|\mathbf{v}\|^2) = \|E(\mathbf{v})\|^2 + \text{tr}\{\text{Cov}(\mathbf{v})\}.$$

**Result 1.11.7.** Let  $\mathbf{v}$  be a random vector and let  $\mathbf{A}$  be a fixed matrix with the same

number of rows as  $\mathbf{v}$ . Then

$$E(\mathbf{v}^\top \mathbf{A} \mathbf{v}) = E(\mathbf{v}) \mathbf{A} E(\mathbf{v})^\top + \text{tr}\{\mathbf{A} \text{Cov}(\mathbf{v})\}.$$

**Result 1.11.8.** Let  $\mathbf{x}$  and  $\mathbf{y}$  be  $n \times 1$  random vectors such that  $\mathbf{x}$  is conditioned on  $\mathbf{y}$ . Then,

$$\text{Cov}(\mathbf{y}) = E\{\text{Cov}(\mathbf{y}|\mathbf{x})\} + \text{Cov}\{E(\mathbf{y}|\mathbf{x})\}.$$

## 1.12 Statistical software

Many longitudinal and multilevel models can be couched in the framework of Bayesian hierarchical models and thus have natural graphical model representations. However, the posterior density functions arise in these types of models are often difficult to obtain analytically. Hence one relies on taking advantage of the growing body of software for approximate inference in general graphical models.

In recent years, graphical model-based Bayesian inference engines have emerged for facilitating such approximate inference for general classes of parametric and nonparametric regression models (e.g. Luts, 2015; Marley and Wand, 2010). Examples include BUGS (Bayesian inference using Gibbs Sampling) (Ligges *et al.*, 2009), Stan (Stan Development Team, 2015), Infer.NET (Minka *et al.*, 2009) and VIBES (Variational Inference for Bayesian Networks)(Bishop *et al.*, 2002). Each engine offers a different form of the approximate inference. BUGS assumes the model is specified in a directed acyclic graph form and uses relatively slow but highly accurate MCMC Gibbs sampling for inference. Stan is a probabilistic programming language written in C++ and uses the Hamiltonian Monte Carlo No U-turn sampling for implementing full Bayesian statistical inference. Both engines are easy to use, especially since they can be directly called from the R computing environment via the wrapper packages BRugs (Ligges *et al.*, 2009) and rstan (Stan Development Team, 2015) respectively. The other two inference engines are designed for performing variational inference using graphical models, but suffer from versatility limitations. Infer.NET uses faster, but less accurate, deterministic approximations such as mean field variational Bayes and expectation propagation. VIBES is open-source Java and is similar to BUGS in terms of functionality. It also uses variational algorithms for inference, but is limited to the conjugate exponential family. Throughout this thesis, we use rstan for MCMC sampling.



## Chapter 2

# Mean Field Variational Bayes Approximations for Longitudinal and Multilevel Data Analysis

*An approximate answer to the right problem is worth a good deal more than an exact answer to an approximate problem.*

John Tukey

### 2.1 Introduction

The past decades have seen the emergence of machine learning and data mining tools designed to help capture and understand patterns from large datasets with complex structure. The most common data types are the longitudinal and multilevel data, which have a grouped or hierarchical structure and are frequently seen in many applied areas such as education, epidemiology, medicine and social science. Examples include human growth studies in which individuals' body measurements are collected at multiple follow-up times, medical studies in which patients are grouped within hospitals and hospital-specific disease rates are examined, and sample surveys in which respondents to questionnaires are grouped within households and geographical districts. These data structures give rise to correlations among observations within groups, therefore sophisticated statistical models are required to analyse these data taking into account correlations induced in the data. A

---

The main content of this chapter is published as: Lee, C. Y. Y. and Wand, M. P. (2015). Streamlined mean field variational Bayes for longitudinal and multilevel data analysis. *Biometrical Journal article in press*. This research has been accepted for presentation at four conferences.

popular model is the *generalised linear mixed model* (GLMM), which directly acknowledges multiple levels of dependency and accounts for different response types.

GLMMs extend generalised linear models by adding normal random effects on the linear predictor scale, to account for correlated observations within groups. However, this flexibility is traded with analytical tractability and may lead to negative implications on computational complexity and efficiency. In the frequentist paradigm, estimation of GLMMs using maximum likelihood is challenging since the integral over the random effects is intractable. Various approximate methods to such problems have been developed in the past two decades, including a seminal approach first introduced by Breslow and Clayton (1993), namely penalised quasi-likelihood, Laplace approximation and its extension (e.g. Raudenbush *et al.*, 2000), and Gaussian variational approximation (e.g. Ormerod and Wand, 2012). In the Bayesian paradigm, an early work by Zeger and Karim (1991) describe approximate Gibbs sampling for GLMMs, using normal distributions to approximate the nonstandard conditional distributions. Fong *et al.* (2009) propose the integrated nested Laplace approximation as a computationally convenient alternative to MCMC. Nonetheless, many of the aforementioned methods involving numerical quadrature or MCMC techniques to approximate intractable integrals suffer from computational intensity.

In this chapter, we examine the use of variational approximations, in particular MFVB, for fitting and inference in Bayesian GLMMs, with an emphasis on the Gaussian, Student- $t$ , Bernoulli and Poisson responses. Variational methods (e.g. Bishop, 2006; Ormerod and Wand, 2010) have evolved mainly in areas such as machine learning and pattern recognition. However, since the early 2000s there has been an increasing recognition of their usefulness in mainstream statistics (e.g. Titterton, 2004). An early contribution is Teschendorff *et al.* (2005), who apply MFVB to mixture modelling for gene expression data. Longitudinal and multilevel models (Diggle *et al.*, 2002; Fitzmaurice *et al.*, 2012; Gelman and Hill, 2007; Goldstein, 2011) represent a major branch of statistics which, to date, has had a relatively little overlap with the graphical models viewpoint and variational methods. Recall from (1.7) that MFVB approximations arise upon restricting  $q(\boldsymbol{\theta})$  to some class of density functions as follows:

$$\mathcal{Q} = \left\{ q(\boldsymbol{\theta}) : q(\boldsymbol{\theta}) = \prod_{i=1}^M q_i(\boldsymbol{\theta}_i) \text{ for some partition } \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\} \text{ of } \boldsymbol{\theta} \right\}.$$

The class of density functions  $\mathcal{Q}$  can either be parametric, semiparametric or nonparametric. Parametric MFVB involves setting  $\mathcal{Q}$  to be a parametric family of density functions. A recent contribution of this type is Opper and Archambeau (2009) and Challis and Barber (2013) who use the Gaussian families. Semiparametric MFVB involves the relaxation of the ordinary mean-field restriction in which some of the density functions in the postu-

lated product density form are pre-specified to be particular parametric density functions (Rhode and Wand, 2015). Nonparametric MFVB means that there is no parametric specification at all on  $\mathcal{Q}$  and the product restriction is the only pre-specification being made. The parametric and semiparametric approaches to restricting  $q(\boldsymbol{\theta})$  are the focus of this chapter.

This chapter makes three major contributions. First, we show, step by step, how to implement MFVB algorithms for Bayesian GLMMs with four common response distributions and their semiparametric extensions. Second, we empirically assess the performance of MFVB algorithms, benchmarking against the standard MCMC methods, through a series of comprehensive simulations. Finally, we derive a variational lower bound for each model that can be useful for model selection. Ultimately, we demonstrate that variational approximations make feasible inferences and estimation of models that would be difficult or impossible to estimate using standard MCMC methods.

Section 2.2 introduces longitudinal and multilevel models and shares a graph theoretical viewpoint of these models, focusing on random intercepts models, random intercepts and slopes models, and their semiparametric regression extensions. Details on the direct implementation of variational algorithms and inference for the four response models are given in Sections 2.5-2.9. In Section 2.10 we provide numerical evidence of the efficacy of variational methods, in terms of both inferential accuracy, credible interval coverage and computational speed. The concluding remarks are given in Section 2.11.

## 2.2 A brief introduction to longitudinal and multilevel models

Longitudinal and multilevel data are common in many applied areas such as education, medicine, epidemiology and social science. Such data have a complex grouped structure. An example of this is depicted in Figure 2.1. The data consist of mathematics test scores of 728 students from 48 schools in inner London, United Kingdom. Details of the data are described in Goldstein (2010).

In any complex structure we can identify atomic units that are at the lowest level of the system. For multilevel data, these atomic units are often individuals. Individuals are then grouped into *higher-level* units, commonly referred to as “groups” or “clusters”. By convention we then say that the individuals are at Level 1 and the higher-level units are at Level 2 in our structure. For longitudinal data, repeat measurements on different occasions are at Level 1 and the individuals are at Level 2.

When individuals form groups or clusters, we might expect that two randomly selected individuals from the same group will be more alike than two individuals selected from different groups. For example, patients in the same hospital share common risk profiles

## 2.2. A BRIEF INTRODUCTION TO LONGITUDINAL AND MULTILEVEL MODELS

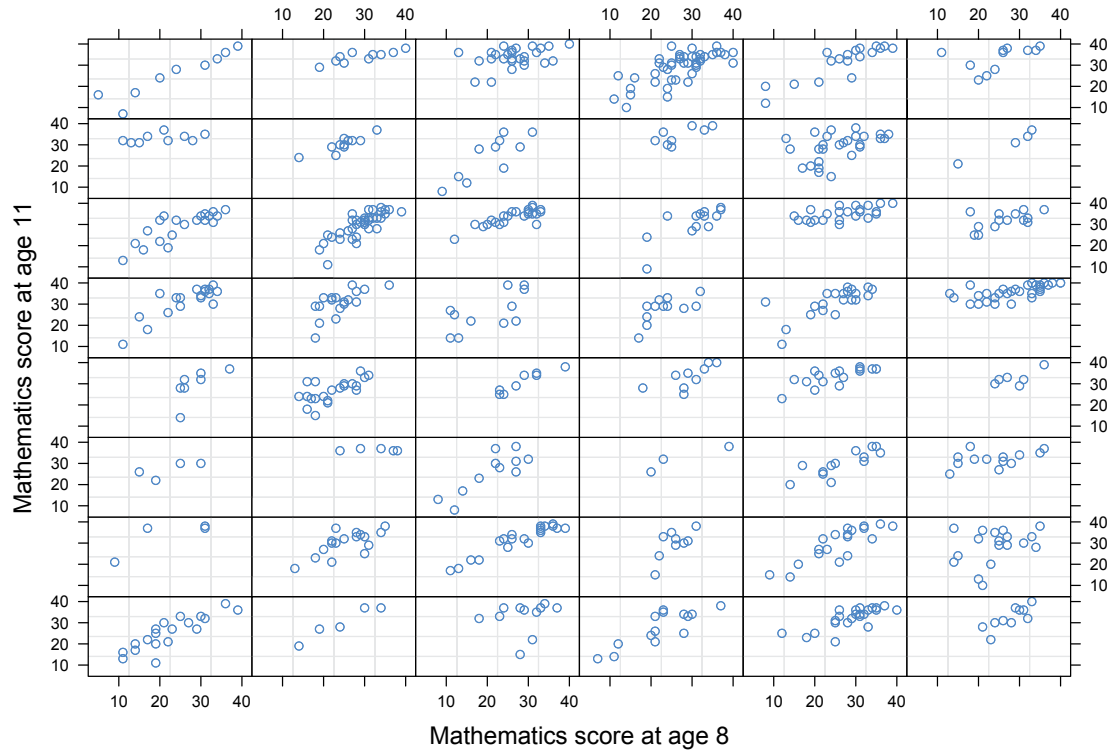


Figure 2.1: Panel of scatterplots of the 11-year old mathematics test scores versus the eight-year old test scores. Each panel corresponds to a school. Details of the data are described in Goldstein (2010).

and hospital characteristics, such as clinical management and practice, are likely to have a similar clinical outcome. By similar argument, measurements taken on the same individual over time will be more highly correlated than two measurements from different individuals. Such correlation violates the independence assumption of the classic regression model and therefore requires a more sophisticated model to appropriately take the induced correlation into account, namely *longitudinal or multilevel models*. Longitudinal or multilevel models are also known as hierarchical models, mixed models and variance components models. For the sake of consistency, we use the term *mixed models* throughout this thesis to represent any general classes of longitudinal and multilevel models, unless stated otherwise. In addition, we restrict our attention to data with a *nested hierarchical* structure, where lower-level unit nests in one and only one higher-level unit. However, it would be relatively straightforward to extend our methods below to include non-nested random effects models.

### 2.2.1 Graph theoretical viewpoint of longitudinal and multilevel models

Graphical models have become increasingly viewed as a general Bayesian inference engine. In many instances, sophisticated probabilistic models are embedded into the framework of graphical models in order to simplify statistical inference and improve computational efficiency. Graphical models are a great way to represent conditional independence assumptions by using graphs. Specifically, nodes represent random variables and lack of edges represent conditional independencies. In addition, the graph is a useful visual representation for complex stochastic systems. The graphical structure is also the basis of efficient inference algorithms.

### 2.2.2 Random intercepts models

Figure 2.1 is a frequently cited example in education, where students are grouped or clustered within schools. A natural starting point to analysing such a data structure is to assume that the overall mean and deviation of the  $i$ th school from that overall mean are simply straight lines. This leads to the standard random intercept model (Laird and Ware, 1982). Consider the Bayesian version of a two-level Gaussian mixed model with a random intercept for each school, and with simplicity in mind, we confine the model to the case of a single predictor, say  $\mathbf{x}$ ,

$$\begin{aligned}
 y_{ij} | \beta_0, \beta_x, u_i^R, \sigma_\varepsilon^2 &\stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_x x_{ij} + u_i^R, \sigma_\varepsilon^2) & (2.1) \\
 \beta_0 &\sim N(\mu_{\beta_0}, \sigma_{\beta_0}^2), \quad \beta_x \sim N(\mu_{\beta_x}, \sigma_{\beta_x}^2), \\
 u_i^R | \sigma_R^2 &\stackrel{\text{ind.}}{\sim} N(0, \sigma_R^2), \\
 \sigma_R^2 &\sim \text{Inverse-Gamma}(A_R, B_R) \\
 \text{and } \sigma_\varepsilon^2 &\sim \text{Inverse-Gamma}(A, B), \quad 1 \leq j \leq n_i, 1 \leq i \leq m.
 \end{aligned}$$

The letters  $i, j$  denote the *nested* indexing for the group level (Level 2) and unit level (Level 1) respectively. Further,  $y_{ij}$  denotes the  $j$ th response at Level 1 within the  $i$ th group at Level 2,  $x_{ij}$  denotes the  $j$ th predictor within the  $i$ th group,  $\beta_0$  and  $\beta_x$  are the so-called *fixed* intercept and slope respectively, and  $u_i^R$  is the so-called *random* intercept for each group. The *mixed model* terminology is due to the fact that model (2.1) is a statistical model containing both *fixed effects* and *random effects*. Here we use the random intercepts to model the dependence of the  $y$ s within higher levels. It is also reasonable to assume that the level parameters are independent of each other and each of them is independent of the error terms  $\varepsilon_{ij}$ .

Recall from Section 1.4 that DAGs provide a useful “road map” of the model’s structure and aid the algebra required for variational approximations. A DAG representation of model (2.1) is depicted in Figure 2.2. In graphical model terminology,  $\mathcal{E} = \{y_{11}, \dots, y_{mn_m}\}$

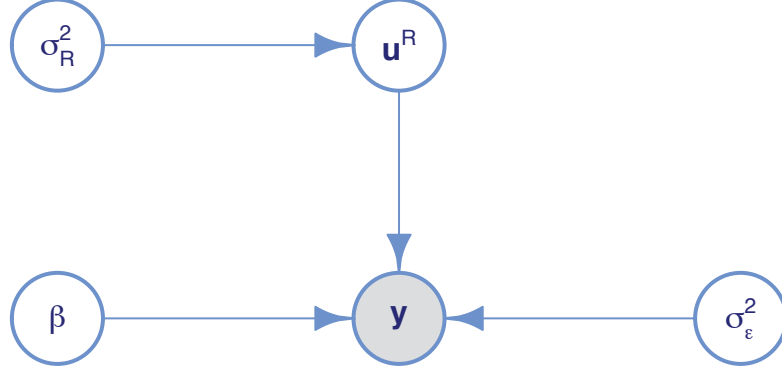


Figure 2.2: Directed acyclic graph corresponding to model (2.1). The open circles represent the random nodes. The shaded circle represents the observed node. The vector  $\mathbf{y} = (y_{11}, \dots, y_{mn_m})$  is represented by a single node, and  $\beta_0$  and  $\beta_x$  are represented by the  $\beta$  node. Similarly,  $\mathbf{u}^R = (u_1^R, \dots, u_m^R)$ .

is the evidence node and  $\mathcal{H} = \{\beta_0, \beta_x, u_1^R, \dots, u_m^R, \sigma_{\beta_0}^2, \sigma_{\beta_x}^2, \sigma_R^2, \sigma_\varepsilon^2\}$  is the set of hidden nodes. Bayesian inference relies upon

$$\begin{aligned} p(\mathcal{H}|\mathcal{E}) &= p(\beta_0, \beta_x, u_1^R, \dots, u_m^R, \sigma_{\beta_0}^2, \sigma_{\beta_x}^2, \sigma_R^2, \sigma_\varepsilon^2 | y_{11}, \dots, y_{mn_m}) \\ &= \frac{p(\mathcal{H}, \mathcal{E})}{\mathcal{E}} = \frac{p(\beta_0, \beta_x, u_1^R, \dots, u_m^R, \sigma_{\beta_0}^2, \sigma_{\beta_x}^2, \sigma_R^2, \sigma_\varepsilon^2, y_{11}, \dots, y_{mn_m})}{p(y_{11}, \dots, y_{mn_m})}, \end{aligned}$$

the posterior density function of all of the random nodes given the observed node. Quite often we are stuck with intractable integrals arising from integrating over the variance components, and therefore obtaining the posterior density function via direct calculation becomes infeasible.

### 2.2.3 Random intercepts and slopes models

Model (2.1) can be easily extended to include both a random intercept and random slope for each school as follows:

$$\begin{aligned} y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2 &\stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_x x_{ij} + u_{0i}^R + u_{1i}^R x_{ij}, \sigma_\varepsilon^2), \\ \beta_0 &\sim N(\mu_{\beta_0}, \sigma_{\beta_0}^2), \quad \beta_x \sim N(\mu_{\beta_x}, \sigma_{\beta_x}^2), \\ [u_{0i}^R \quad u_{1i}^R]^\top | \Sigma^R &\stackrel{\text{ind.}}{\sim} N\left(\mathbf{0}, \Sigma^R \equiv \begin{bmatrix} \sigma_{R_0}^2 & \sigma_{R_0, R_1} \\ \sigma_{R_0, R_1} & \sigma_{R_1}^2 \end{bmatrix}\right), \\ \Sigma^R &\sim \text{Inverse-Wishart}(A_\Sigma, \mathbf{B}_\Sigma) \\ \text{and } \sigma_\varepsilon^2 &\sim \text{Inverse-Gamma}(A, B), \quad 1 \leq j \leq n_i, 1 \leq i \leq m. \end{aligned} \tag{2.2}$$

### 2.3. MIXED MODEL REPRESENTATION

---

The majority of the parameters are defined analogously in Subsection 2.2.2. The terms  $u_{0i}^R$  and  $u_{1i}^R$  are the so-called random intercepts and random slopes, being treated as a random sample from a bivariate normal distribution with an unstructured  $2 \times 2$  covariance matrix  $\Sigma^R$ . This model accounts for possible variability in the group intercepts  $\sigma_{R_0}^2$  and slopes  $\sigma_{R_1}^2$  and allows for an intercept-slope correlation  $\rho_{R_0, R_1}$ .

#### 2.2.4 Extension to semiparametric regression

Semiparametric regression models extend generalised linear mixed models to accommodate non-linear predictor effects. The essence of semiparametric regression is to include penalised spline basis functions through random effects formulation that have the same form as those used traditionally in longitudinal and multilevel data analysis (e.g. Ruppert *et al.*, 2003).

Here we extend model (2.2) by simply adding a smooth, but otherwise unspecified, function for the  $\ell$ th predictor, say  $s_\ell$ ,  $1 \leq \ell \leq L$ ,

$$f_\ell(s_\ell) = \beta_\ell s_\ell + \sum_{k=1}^{q_\ell^G} u_{\ell k}^G z_{\ell k}(s_\ell), \quad 1 \leq \ell \leq L, \quad (2.3)$$

where  $z_{\ell 1}(\cdot), \dots, z_{\ell q_\ell^G}(\cdot)$  is a set of spline functions appropriate for the linear component being unpenalised. Our default choice for the  $z_\ell(\cdot)$  is the B-spline basis and penalty set-up of O'Sullivan (1986), since this leads to approximate smoothing splines which have good boundary and extrapolation properties. The quantity of spline basis functions has a minimal effect on the adequacy of (2.3) and  $q_\ell^G = 25$  is recommended for models of practical interest (Li and Ruppert, 2008). The coefficients  $u_{\ell k}^G$  can be considered as a measure of the basis amplitude since they regulate the contribution from the  $\ell$ th curve.

From a computational point of view, penalised splines are equivalent to the random effects, in that we penalise the spline basis function coefficients by treating them as a random sample from a multivariate normal distribution to avoid overfitting of the data. That is,

$$u_{11}^G, \dots, u_{Lq_L^G}^G | \sigma_{u1}^2, \dots, \sigma_{uL}^2 \sim N \left( \mathbf{0}, \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} \right).$$

## 2.3 Mixed model representation

Classical mixed models can be viewed as an intermediate step between frequentist and Bayesian models since the grouped effects are treated as random. Bayesian analysis treats all model parameters as random, assigns prior distributions to characterise our prior knowledge about these parameters prior to data collection, and computes the joint posterior

### 2.3. MIXED MODEL REPRESENTATION

---

density function of all of the random parameters given the observed data as the basis of inference.

In common practice, we generalise the fixed and random components of the mixed models to arbitrary general design matrices (Zhao *et al.*, 2006). This allows one to take advantage of the ever-expanding methods and software for inference in these models. Henceforth, we combine models (2.2) and (2.3) in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad (2.4)$$

where

$$\begin{aligned} \mathbf{y} &\equiv \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X}_i \equiv \begin{bmatrix} \mathbf{1} & \mathbf{x}_i \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} \mathbf{X}_1 & \mathbf{s}_1 \\ \vdots & \vdots \\ \mathbf{X}_m & \mathbf{s}_m \end{bmatrix}, \quad \boldsymbol{\varepsilon} \equiv \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \vdots \\ \boldsymbol{\varepsilon}_m \end{bmatrix}, \\ \mathbf{Z} &\equiv \begin{bmatrix} z_{11}(\mathbf{s}_{11}) & \dots & z_{1q_1^G}(\mathbf{s}_{11}) & \dots & z_{L1}(\mathbf{s}_{L1}) & \dots & z_{Lq_L^G}(\mathbf{s}_{L1}) & \mathbf{X}_1 & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots & \dots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ z_{11}(\mathbf{s}_{1m}) & \dots & z_{1q_1^G}(\mathbf{s}_{1m}) & \dots & z_{L1}(\mathbf{s}_{Lm}) & \dots & z_{Lq_L^G}(\mathbf{s}_{Lm}) & \mathbf{0} & \dots & \mathbf{X}_m \end{bmatrix}, \\ \boldsymbol{\beta} &\equiv [\beta_0 \quad \beta_x \quad \beta_{s1} \quad \dots \quad \beta_{sL}]^\top \quad \text{and} \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}^G \\ \mathbf{u}_1^R \\ \vdots \\ \mathbf{u}_m^R \end{bmatrix}. \end{aligned}$$

For the  $i$ th group,  $\mathbf{y}_i$  is an  $n_i \times 1$  vector of response variables,  $\mathbf{x}_i$  and  $\mathbf{s}_i$  are  $n_i \times 1$  vectors of predictor variables,  $\boldsymbol{\varepsilon}_i$  is an  $n_i \times 1$  vector of errors and  $\mathbf{u}_i^R$  is a  $q^R \times 1$  vector comprising the random intercept and random slope of the  $i$ th group. The  $\sum_{\ell=1}^L q_\ell^G \times 1$  vector of spline basis coefficients is denoted by  $\mathbf{u}^G \equiv [\mathbf{u}_1^G \dots \mathbf{u}_L^G]$ . The matrices  $\mathbf{X}$  and  $\mathbf{Z}$  are the respective  $\sum_{i=1}^m n_i \times p$  fixed effects design matrix and  $\sum_{i=1}^m n_i \times q$  random effects design matrix,  $\boldsymbol{\beta}$  and  $\mathbf{u}$  are the  $p \times 1$  fixed effects vector and  $q \times 1$  random effects vector. The random effects have a multivariate normal distribution with zero mean and covariance matrix  $\mathbf{G}$ , given as

$$\mathbf{G} = \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^R) \end{bmatrix}.$$

Putting all of these components together gives the following *linear mixed model* form

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), \quad \mathbf{u} | \mathbf{G} \sim N(\mathbf{0}, \mathbf{G}). \quad (2.5)$$

A common extension of (2.5) involves replacement of  $\sigma_\varepsilon^2 \mathbf{I}$  by a general covariance matrix  $\mathbf{R}$ , but is not treated here.



A very wide range of models can be obtained through various choices of  $\mathbf{Z}$  and  $\mathbf{G}$ . In this chapter, we focus on those corresponding to two-level hierarchical structures, which are often treated separately within the broad branches of longitudinal data analysis (Diggle *et al.*, 2002; Fitzmaurice *et al.*, 2012) and multilevel modelling (Gelman and Hill, 2007; Goldstein, 2011). Since around 2000, there has been a major interplay between longitudinal data analysis and semiparametric regression. An overview is given in Fitzmaurice *et al.* (2008). As we shall see in subsequent chapters, the general model (2.4) is extremely rich in that it encompasses a wide range of models that are useful in practice.

The class of generalised linear mixed models is much more general than (2.5), as explained by Zhao *et al.* (2006). Staying within the one-parameter exponential family, a more general class of models is

$$\mathbf{y}|\boldsymbol{\beta}, \mathbf{u} \sim \exp \left\{ \mathbf{y}^\top (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top b(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) + \mathbf{1}^\top c(\mathbf{y}) \right\}, \quad \mathbf{u}|\mathbf{G} \sim N(\mathbf{0}, \mathbf{G}), \quad (2.6)$$

where the design matrices  $\mathbf{X}$  and  $\mathbf{Z}$  and covariance matrix  $\mathbf{G}$  are quite general and described above. The functions  $b$  and  $c$  are specific to members of the family. The most common examples are Bernoulli for which  $b(x) \equiv \log(1 + e^x)$  and  $c(x) \equiv 0$  and Poisson for which  $b(x) \equiv e^x$  and  $c(x) \equiv -\log(x!)$ .

### 2.3.1 Sample size and subscript notation

Many of the models used in longitudinal and multilevel data analysis are mathematically equivalent, but use different sample size and subscript notation. It is prudent to delineate the various notational conventions commonly in use so that methodology developed using one set of notation can be transferred to models that use other notational systems. We achieve that in this Subsection by way of a concrete example of (2.5).

The special case of (2.5) with

$$\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad \mathbf{Z} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{G} = \sigma_{\text{R}}^2 \mathbf{I} \quad (2.7)$$

### 2.3. MIXED MODEL REPRESENTATION

---

imposes the following covariance matrix on the response vector:

$$\text{Cov}(\mathbf{y}) = \begin{bmatrix} \sigma_R^2 + \sigma_\varepsilon^2 & \sigma_R^2 & 0 & 0 & 0 & 0 & 0 \\ \sigma_R^2 & \sigma_R^2 + \sigma_\varepsilon^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_R^2 + \sigma_\varepsilon^2 & \sigma_R^2 & \sigma_R^2 & 0 & 0 \\ 0 & 0 & \sigma_R^2 & \sigma_R^2 + \sigma_\varepsilon^2 & \sigma_R^2 & 0 & 0 \\ 0 & 0 & \sigma_R^2 & \sigma_R^2 & \sigma_R^2 + \sigma_\varepsilon^2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sigma_R^2 + \sigma_\varepsilon^2 & \sigma_R^2 \\ 0 & 0 & 0 & 0 & 0 & \sigma_R^2 & \sigma_R^2 + \sigma_\varepsilon^2 \end{bmatrix}.$$

This block covariance structure is in keeping with two hierarchical levels with 3 groups at the first level and sample sizes of 2, 3 and 2 observations within each group.

There are at least two established notations for conveying the hierarchical structure imposed by (2.7). The first is

$$E(y_{ij} | u_i^R) = u_i^R,$$

where

$$\begin{aligned} y_{ij} &\equiv j\text{th measurement in the } i\text{th group,} \\ u_i^R &\equiv \text{random effect in the } i\text{th group,} \\ &\text{for } 1 \leq i \leq m, \quad 1 \leq j \leq n_i. \end{aligned} \tag{2.8}$$

In the current example,  $m = 3$  and  $n_1 = 2, n_2 = 3, n_3 = 2$ . This notation is used by prominent longitudinal data analysis textbooks, Diggle *et al.* (2002) and Fitzmaurice *et al.* (2012), as well as the semiparametric regression book, Ruppert *et al.* (2003).

An alternative notation is

$$E(y_i | u_{j[i]}^R) = u_{j[i]}^R, \quad 1 \leq i \leq 7, \quad 1 \leq j \leq J,$$

where for this example,  $J = 3$  and the  $j[i]$  mapping is

$$1[1] \equiv 1, \quad 1[2] \equiv 1, \quad 1[3] \equiv 1, \quad 2[3] \equiv 2, \quad 2[4] \equiv 2, \quad 2[5] \equiv 2, \quad 3[6] \equiv 3, \quad 3[7] \equiv 3$$

and

$$y_i = i\text{th entry of } \mathbf{y}.$$

This is used in the multilevel model textbook by Gelman and Hill (2007). Goldstein (2011) uses

$$y_{ij} \equiv i\text{th measurement in the } j\text{th group,} \quad 1 \leq j \leq m$$

for two-level models with  $m$  groups.

For the remainder of this thesis we will use the sample size and subscript notation adopted by the above-mentioned longitudinal data analysis textbooks. This involves sub-

script notation corresponding to (2.8) and

$m \equiv$  number of groups and  $n_i \equiv$  number of response measurements in the  $i$ th group.

The mathematical equivalence between longitudinal and multilevel models means that our methodology applies equally to both areas. However, as illustrated above, some notational adjustment is required to match common multilevel model notations.

## 2.4 Prior distributions for fixed and random effects

Throughout this thesis we take the prior distribution for the fixed effects vector  $\boldsymbol{\beta}$ , standard deviation parameters for the spline bases  $\sigma_{ul}$  and random effects covariance matrix  $\boldsymbol{\Sigma}^R$  of the form:

$$\begin{aligned} \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}_p), \quad \sigma_{ul} \overset{\text{ind.}}{\sim} \text{Half-Cauchy}(A_{ul}) \\ \text{and} \quad \boldsymbol{\Sigma}^R &\sim \text{Inverse-Wishart}(A_{\boldsymbol{\Sigma}}, \mathbf{B}_{\boldsymbol{\Sigma}}), \end{aligned} \quad (2.9)$$

where the hyperparameters  $\sigma_{\boldsymbol{\beta}}^2$ ,  $A_{ul} > 0$  are to be specified by the user. The standard deviation parameters have independent half-Cauchy priors. As explained in Gelman (2006), (2.9) is an attractive means by which weakly-informative priors can be imposed on the  $\sigma_{ul}$ s. Result 5 of Wand *et al.* (2012) leads to the following equivalent distributional statement:

$$\sigma_{ul}^2 \overset{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{ul}\right), \quad a_{ul} \overset{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{ul}^2\right). \quad (2.10)$$

We use representation (2.10) rather than (2.9) since it is more conducive to MFVB.

For  $\boldsymbol{\Sigma}^R$  we use the following extension of (2.10) (Huang and Wand, 2013):

$$\begin{aligned} \boldsymbol{\Sigma}^R | a_1^R, \dots, a_{q^R}^R &\sim \text{Inverse-Wishart}\left(\nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R)\right), \\ a_r^R &\overset{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{Rr}^2\right), \quad r = 1, \dots, q^R, \end{aligned}$$

where  $\nu, A_{Rr} > 0$ . The choice  $\nu = 2$  corresponds to the correlation parameters having uniform distributions over  $(-1, 1)$  and the standard deviation parameters having Half- $t$  distributions with 2 degrees of freedom.

For notation convenience, define the following parameter and matrices throughout this chapter as follows:

$$N \equiv \sum_{i=1}^m n_i, \quad \mathbf{C} \equiv [\mathbf{X} \quad \mathbf{Z}], \quad \mathbf{C}_{12} \equiv \text{diag}(b_i)_{1 \leq i \leq N} \quad \text{and} \quad \mathbf{v} \equiv \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}.$$

The effects vector  $\mathbf{v}$  has a multivariate normal distribution with zero mean and covariance

matrix  $\mathbf{\Omega}$ , given as

$$\mathbf{v} | \sigma_\beta^2, \sigma_{u1}^2, \dots, \sigma_{uL}^2, \mathbf{\Sigma}^R \sim N \left( \mathbf{0}, \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\sigma_{u\ell}^{-2} \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes (\mathbf{\Sigma}^R)^{-1} \end{bmatrix} \right).$$

In what follows, we systematically present approximate Bayesian analysis of four common GLMMs in scientific research.

## 2.5 Gaussian semiparametric mixed models

We first consider the case where the distribution of response vector  $\mathbf{y}$  is Gaussian, the two-level Bayesian Gaussian semiparametric mixed model stated in full is

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), & \boldsymbol{\beta} &\sim N(0, \sigma_\beta^2 \mathbf{I}_p), & (2.11) \\ \mathbf{u} | \sigma_{u1}^2, \dots, \sigma_{uL}^2, \mathbf{\Sigma}^R &\sim N \left( \mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \mathbf{\Sigma}^R \end{bmatrix} \right), \\ \mathbf{\Sigma}^R | a_1^R, \dots, a_{q^R}^R &\sim \text{Inverse-Wishart} \left( \nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R) \right), \\ a_r^R &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{Rr}^2 \right), & r &= 1, \dots, q^R, \\ \sigma_{u\ell}^2 | a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/a_{u\ell} \right), & a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{u\ell}^2 \right), \\ \sigma_\varepsilon^2 | a_\varepsilon &\sim \text{Inverse-Gamma} \left( \frac{1}{2}, 1/a_\varepsilon \right), & a_\varepsilon &\sim \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_\varepsilon^2 \right). \end{aligned}$$

Figure 2.3 shows the graph-theoretic representation of model (2.11). The advantages to such a graphical representation are two-fold: first, it provides a visualisation of the hierarchical structure of its corresponding Bayesian model; and second, graph-theoretic results can be used to determine probabilistic relationships between nodes. The node  $\mathbf{a}^R$  corresponds to the random vector  $[a_1^R \dots a_{q^R}^R]^\top$ . The nodes  $\boldsymbol{\sigma}_u^2$  and  $\mathbf{a}_u$  are defined in a similar vein, i.e.  $\boldsymbol{\sigma}_u^2 \equiv [\sigma_{u1}^2 \dots \sigma_{uL}^2]^\top$  and  $\mathbf{a}_u \equiv [a_{u1}^2, \dots, a_{uL}^2]^\top$  respectively. The node  $\mathbf{u}$  is separated into two nodes  $\mathbf{u}^R$  and  $\mathbf{u}^G$ . Note that the random effects have been partitioned into spline coefficients (superscript G) and group effects (superscript R).

### 2.5.1 Approximate Bayesian inference via mean field variational Bayes

Consider the problem of Bayesian inference for the random parameters in model (2.11) as depicted in Figure 2.3. The process of MFVB approximation is best illustrated in Figure 2.5 using “life cycle of nodes” as a metaphor. We begin with the DAG that can be viewed as a family tree, where each directed edge signifies a “parent-child” relationship. Moralisation then follows, which involves “marrying” any unconnected parents of a common child. The resultant relationships are complex and often lead to an intractable problem. This is where

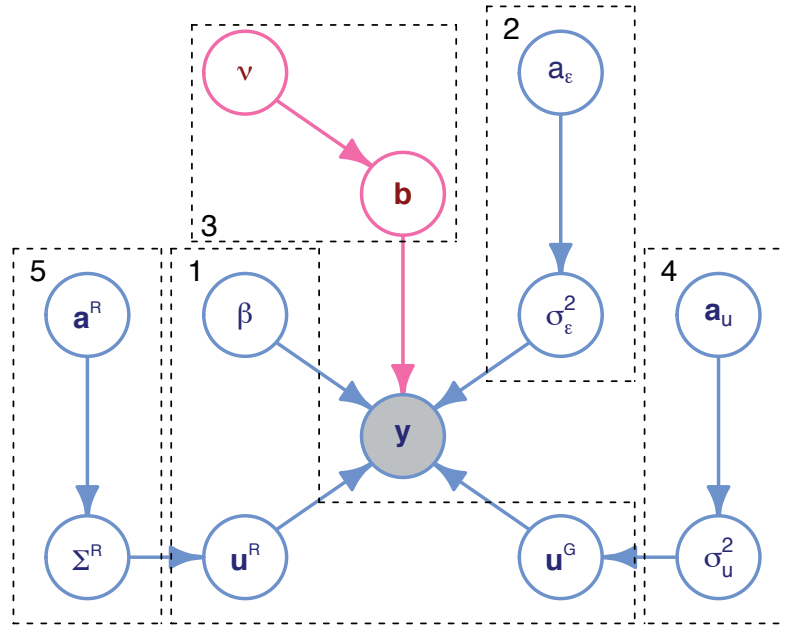


Figure 2.3: Directed acyclic graph for the two-level Bayesian semiparametric mixed models with the Gaussian and Student- $t$  responses. The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The pink open nodes  $b$  and  $\nu$  are additional parameters for the Student- $t$  response model. The numbered grey dashed polygons indicate the corresponding update expressions in Algorithm 3.

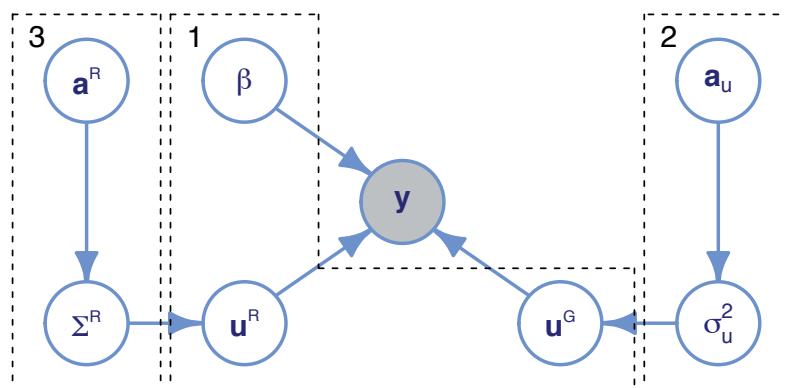


Figure 2.4: Directed acyclic graph for the two-level Bayesian semiparametric mixed models with the Bernoulli and Poisson responses. The shaded node corresponds to the observed data vector and the open node corresponds to the random or auxiliary variables. The numbered grey dashed polygons indicate the corresponding update expressions in Algorithms 4 and 6.

the MFVB approximation comes into play by “getting rid of some kids” and assumes posterior independence among parameters in order to achieve tractability.

To elaborate mathematically, the essence of this approach involves approximation of the full joint posterior density function of the form

$$p(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2 | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2)$$

subject to the  $q$ -density product restriction

$$q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2) = q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon) \times q(\boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2). \quad (2.12)$$

Restriction (2.12) is the minimal factorisation in the MFVB approximation. Menictas & Wand (2013) provides some heuristic arguments for the restriction in (2.12). Essentially, such a restriction is underpinned by the well-known parameter orthogonality result arising from the block-diagonal form of the Fisher information matrix, which implies that the maximum likelihood estimators of  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$  in linear regression models are asymptotically independent. Applying this result to the approximate Bayesian inference context leads to (2.12) and so we would expect this posterior independence assumption to be reasonable, and therefore, should incur little loss in accuracy.

The  $q$ -densities are chosen to minimise the Kullback-Leibler divergence between the full joint posterior density function and the right hand side of (2.12) (Wainwright and Jordan, 2008). Although the mild restriction (2.12) is sufficient to produce closed-form  $q$ -densities, *induced product results* (Bishop, 2006) lead to further factorisations in the product  $q$ -density form that simplify calculations. Such induced factorisations arise from an interaction between the factorisation assumed in the approximating posterior density function and the conditional independence properties of the actual joint density function, and can be easily detected using simple graphical tests such as moralisation. As an illustration, Figure 2.5b is the moral graph of model (2.11). Let  $A = \sigma_\varepsilon^2$ ,  $B = \boldsymbol{\Sigma}^R$  and  $C = \{\boldsymbol{\beta}, \mathbf{u}, \mathbf{y}\}$ , then the smallest ancestral set containing  $A \cup B \cup C$  is the entire DAG. It is evident that all paths connecting nodes in  $A$  and  $B$  must pass through at least one of the nodes in  $C$ . In other words, the set of nodes in  $C$  blocks the paths from  $A$  to  $B$  and therefore the following conditional independence statement holds under the global Markov property:

$$A \perp\!\!\!\perp B \mid C, \quad \text{i.e.} \quad \sigma_\varepsilon^2 \perp\!\!\!\perp \boldsymbol{\Sigma}^R \mid \{\boldsymbol{\beta}, \mathbf{u}, \mathbf{y}\}.$$

Similarly, we can deduce from the global Markov property the conditional independences

$$\{\boldsymbol{\beta}, \mathbf{u}\} \perp\!\!\!\perp \mathbf{a}^R \mid \{\mathbf{y}, \boldsymbol{\Sigma}^R, \sigma_\varepsilon^2\} \quad \text{and} \quad \mathbf{a}^R \perp\!\!\!\perp \mathbf{a}_u \mid \{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}^R, \sigma_\varepsilon^2, \sigma_u^2\}.$$

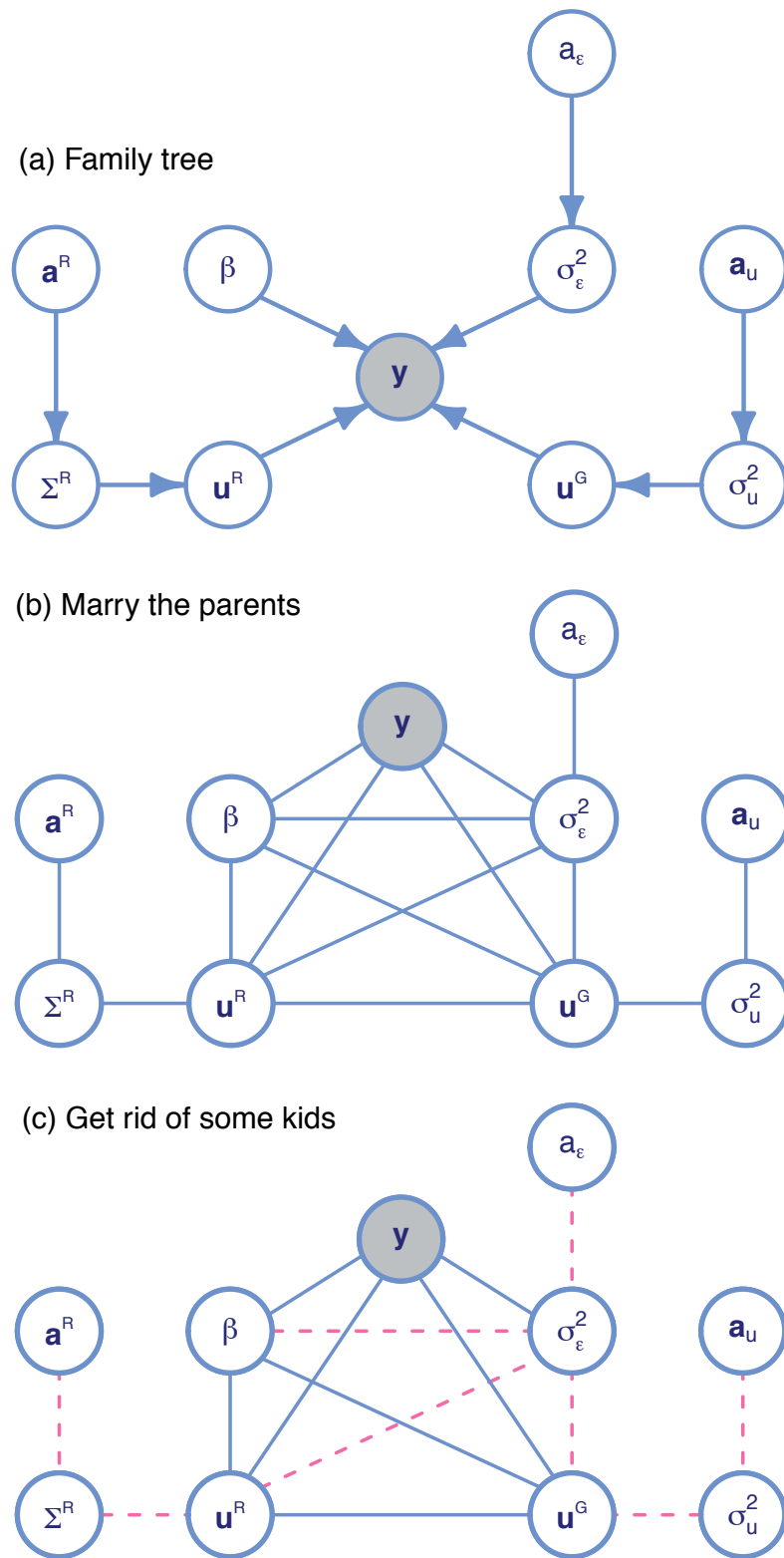


Figure 2.5: Life cycle of nodes for mean field variational Bayes: (a) directed acyclic graph for model (2.11), (b) moral graph for model (2.11), (c) modification of Figure 2.5b with eight edges removed to impose  $q$ -density product restriction. In each graph, shading is used to signify the observed data vector.

Similar arguments lead to the  $q$ -densities having the product form:

$$\begin{aligned} q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2) & \quad (2.13) \\ = q(\boldsymbol{\beta}, \mathbf{u}) q(a_\varepsilon) q(\boldsymbol{\Sigma}^R) q(\sigma_\varepsilon^2) & \left\{ \prod_{r=1}^{q^R} q(a_r^R) \right\} \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\}. \end{aligned}$$

Visually, this is analogous to removing edges between relevant nodes in the moralised DAG, represented by the pink dotted lines in Figure 2.5c. Under the product restriction (2.13), all parameters of model (2.11) have closed-form expressions for their full conditional density function:

$$\begin{aligned} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} & \sim N \left( (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{y}, (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \right), \\ \sigma_\varepsilon^2 \Big|_{\text{rest}} & \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), \frac{1}{2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 + a_\varepsilon^{-1} \right), \\ a_\varepsilon \Big|_{\text{rest}} & \sim \text{Inverse-Gamma} \left( 1, \sigma_\varepsilon^{-2} + A_\varepsilon^{-2} \right), \\ \sigma_{u\ell}^2 \Big|_{\text{rest}} & \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2} (q_\ell^G + 1), \frac{1}{2} \|\mathbf{u}_\ell^G\|^2 + a_\ell^{-1} \right), \\ a_{u\ell} \Big|_{\text{rest}} & \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( 1, \sigma_{u\ell}^{-2} + A_{u\ell}^{-2} \right), \\ \boldsymbol{\Sigma}^R \Big|_{\text{rest}} & \sim \text{Inverse-Wishart} \left( \nu + m + 1, \sum_{i=1}^m \{u_i^R (u_i^R)^\top\} + 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R) \right), \\ a_r^R \Big|_{\text{rest}} & \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2} (\nu + q^R), \nu (\boldsymbol{\Sigma}^R)_{rr}^{-1} + A_{Rr}^{-2} \right), \quad r = 1, \dots, q^R. \end{aligned}$$

Here “rest” denotes the set of other random variables in model (2.11) and  $(\boldsymbol{\Sigma}^R)_{rr}^{-1}$  denotes the  $(r, r)$  entry of  $(\boldsymbol{\Sigma}^R)^{-1}$ . The full conditionals are of standard form, so MCMC methods can be used to construct a Markov chain that targets from the full conditionals that leads to samples from their joint posterior. However, we demonstrate fitting via MFVB approximation since it provides much faster computation than MCMC.

Recall from Section 1.5.3.1 that the optimal  $q$ -density for a parameter vector  $\boldsymbol{\theta}$ , denoted by  $q^*(\boldsymbol{\theta})$ , were shown to satisfy:

$$q_i^*(\boldsymbol{\theta}_i) \propto \exp\{E_{q(-\boldsymbol{\theta}_i)} \log p(\boldsymbol{\theta}_i | \text{Markov blanket of } \boldsymbol{\theta}_i)\},$$

where  $E_{q(-\boldsymbol{\theta}_i)}$  denotes expectation with respect to the  $q$ -densities of all parameters except  $\boldsymbol{\theta}_i$ . As shown through the derivations in Appendix 2.A, the optimal  $q$ -densities admit the following forms:

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}) & \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function,} & (2.14) \\ q^*(\sigma_\varepsilon^2) & \text{ is the Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)} \right) \text{ density function,} \\ q^*(a_\varepsilon) & \text{ is the Inverse-Gamma}(1, B_{q(a_\varepsilon)}) \text{ density function,} \end{aligned}$$



$q^*(\sigma_{ul}^2)$  is the Inverse-Gamma  $\left(\frac{1}{2}(q_\ell^G + 1), B_{q(\sigma_{ul}^2)}\right)$  density function,  
 $q^*(a_{ul})$  is the Inverse-Gamma  $(1, B_{q(a_{ul})})$  density function,  
 $q^*(\Sigma^R)$  is the Inverse-Wishart  $\left(\nu + m + q^R + 1, \mathbf{B}_{q(\Sigma^R)}\right)$  density function,  
 $q^*(a_r^R)$  is the Inverse-Gamma  $\left(\frac{1}{2}(\nu + q^R), B_{q(a_r^R)}\right)$  density function,

for parameters  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , the mean vector and covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $B_{q(\sigma_\varepsilon^2)}$ , the scale parameter of  $q^*(\sigma_\varepsilon^2)$ ,  $B_{q(a_\varepsilon)}$ , the scale parameter of  $q^*(a_\varepsilon)$ ,  $B_{q(\sigma_{ul}^2)}$ , the scale parameter of  $q^*(\sigma_{ul}^2)$ ,  $B_{q(a_{ul})}$ , the scale parameter of  $q^*(a_{ul})$ ,  $B_{q(a_r^R)}$ , the scale parameter of  $q^*(a_r^R)$  and  $\mathbf{B}_{q(\Sigma^R)}$ , the scale parameter of  $q^*(\Sigma^R)$ . The MFVB approximation involves standard density functions and requires only simple closed-form coordinate ascent updates, see Algorithm 3. Convergence of Algorithm 3 for the Gaussian response model can be monitored using successive values of the lower bound on the marginal log-likelihood:

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^R (\nu + q^R - 1) \log(2\nu) - \frac{1}{2} \log(2\pi) \sum_{i=1}^m n_i - \left(\frac{1}{2} q^R + L + 1\right) \log(\pi) \\
 &\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) \\
 &\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}| - \log(\mathcal{C}_{q^R, \nu + q^R - 1}) + \log(\mathcal{C}_{q^R, \nu + m + q^R - 1}) \\
 &\quad - \frac{1}{2} (\nu + m + q^R - 1) \log |\mathbf{B}_{q(\Sigma^R)}| \\
 &\quad + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^G + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{ul}^2)}) \\
 &\quad + \log \Gamma \left\{ \frac{1}{2} (\sum_{i=1}^m n_i + 1) \right\} - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_\varepsilon^2)}) \\
 &\quad - \sum_{r=1}^{q^R} \log(A_{Rr}) + q^R \log \Gamma \left\{ \frac{1}{2} (\nu + q^R) \right\} \\
 &\quad - \sum_{\ell=1}^L \log(A_{u\ell}) - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\
 &\quad + \sum_{r=1}^{q^R} \nu (\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} \mu_{q(1/a_r^R)} - \frac{1}{2} (\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) \\
 &\quad - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)},
 \end{aligned}$$

where  $\boldsymbol{\mu}_{q(\boldsymbol{\beta})}$  is defined to be the subvector of  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$  and similarly,  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}$  is the sub-matrix of  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$ . The terms  $\mathcal{C}_{q^R, \nu + q^R - 1}$  and  $\mathcal{C}_{q^R, \nu + m + q^R - 1}$  are defined in Table 1.1. Details on the derivation for the optimal  $q$ -densities and lower bound expression are deferred to Appendix 2.A.

### 2.5.2 Approximating density functions of entries of $\Sigma^R$

In Section 2.10, we present the results from a series of comprehensive simulation studies that compare the MFVB methodology against MCMC. We do this comparison on the regression coefficients and the entries of the effects covariance matrix  $\Sigma_{11}^R, \Sigma_{12}^R, \Sigma_{21}^R$  and  $\Sigma_{22}^R$ . The approximating density functions of regression coefficients are straightforward to attain, however, the ones related to the entries of the effects covariance matrix require additional effort. The derivations of the approximating density functions of these entries

are given in this subsection.

### 2.5.2.1 Main diagonal entries

The distributions of the main diagonal entries of  $q^*(\boldsymbol{\Sigma}^R)$  are obtained using the following theorem (Menictas, 2015):

**Theorem 2.5.1.** If  $\mathbf{X}$  is an  $n \times n$  positive definite matrix that has an inverse-Wishart distribution with degrees of freedom  $a$  and scale matrix  $\mathbf{B}$ , then the diagonal entries  $\mathbf{X}_{jj}$  are such that:

$$\mathbf{X}_{jj} \sim \text{Inverse-Gamma}\left(\frac{1}{2}(a - n + 1), \frac{1}{2}\mathbf{B}_{jj}\right), \quad 1 \leq j \leq n.$$

From Algorithm 3 we see that the optimal  $q$ -density of  $\boldsymbol{\Sigma}^R$  has the form

$$\begin{aligned} q^*(\boldsymbol{\Sigma}^R) &\sim \text{Inverse-Wishart}\left(A_{q(\boldsymbol{\Sigma}^R)}, \mathbf{B}_{q(\boldsymbol{\Sigma}^R)}\right), \\ \text{where } A_{q(\boldsymbol{\Sigma}^R)} &= \nu + m + q^R - 1 \\ \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} &= \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}) + 2\nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)}). \end{aligned}$$

Using Theorem 2.5.1, the main diagonal entries take the form

$$q^*(\boldsymbol{\Sigma}_{rr}^R) \sim \text{Inverse-Gamma}\left(\frac{1}{2}(\nu + m), \frac{1}{2}(\mathbf{B}_{q(\boldsymbol{\Sigma}^R)})_{jj}\right), \quad 1 \leq r \leq q^R.$$

### 2.5.2.2 Off-diagonal entries

The distributions of the off diagonal entries of an inverse-Wishart random matrix  $\mathbf{X}$  can be obtained analytically via the fast Fourier transform inversion of the characteristic function (Menictas, 2015):

**Theorem 2.5.2.** If  $\mathbf{X}$  is an  $n \times n$  positive definite matrix that has an inverse-Wishart distribution with degrees of freedom  $a$  and scale matrix  $\mathbf{B}$  then the characteristic function of the off-diagonal entries  $\mathbf{X}_{jj'}$  take the form:

$$E\left(e^{itX_{jj'}}\right) = \exp\left\{-\frac{a}{2} \sum_{\ell=1}^n \log(\lambda_\ell)\right\},$$

where  $\lambda_\ell$  are the eigenvalues of the matrix  $\mathbf{I} - 2i\mathbf{T}\mathbf{B}^{-1}$  and  $\mathbf{T}$  is a matrix with the  $(j, j')$  and  $(j', j)$  entries equal to  $t/2$  and all other entries equal to 0.

Unfortunately, there is no straightforward way of obtaining the above characteristic function. We therefore use a Monte Carlo sampling approach, subject to Monte Carlo errors, to find the distributions of the off diagonal elements in  $q^*(\boldsymbol{\Sigma}^R)$ .

## 2.6 Student- $t$ semiparametric mixed models

We now consider the case where the distribution of response vector  $\mathbf{y}$  is Student- $t$

$$\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \nu, \sigma_\varepsilon^2 \sim t((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}), \sigma_\varepsilon^2, \nu), \quad \nu \sim \text{Uniform}(\nu_{\min}, \nu_{\max}), \quad (2.15)$$

where  $\nu_{\max}, \nu_{\min} > 0$ . The low degrees of freedom value  $\nu > 0$  corresponds to the Student- $t$  distribution with very heavy tails, and it is commonly used to model data containing gross outliers.

The notational infrastructure described in Section 2.5 for the Gaussian response model can be directly transferred to the Student- $t$  response case. The only differences are the response distribution specification and inclusion of additional parameters  $\mathbf{b}$  and  $\nu$ . Using Result 1.8.2, we can rewrite (2.15) and thus the full two-level Bayesian Student- $t$  semiparametric mixed model takes the following form:

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \mathbf{b}, \sigma_\varepsilon^2 &\sim N((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}), \sigma_\varepsilon^2 \mathbf{C}_{12}), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), & (2.16) \\ \mathbf{u}|\sigma_{u_1}^2, \dots, \sigma_{u_L}^2, \boldsymbol{\Sigma}^R &\sim N\left(\mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma}^R \end{bmatrix}\right), \\ \boldsymbol{\Sigma}^R|a_1^R, \dots, a_{q^R}^R &\sim \text{Inverse-Wishart}\left(\nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R)\right), \\ a_r^R &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{Rr}^2\right), \quad r = 1, \dots, q^R, \\ \sigma_{u\ell}^2|a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{u\ell}\right), \quad a_{u\ell} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{u\ell}^2\right), \\ \sigma_\varepsilon^2|a_\varepsilon &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_\varepsilon\right), \quad a_\varepsilon \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_\varepsilon^2\right), \\ b_i|\nu &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{\nu}{2}, \nu/2\right) \quad \text{and} \quad \nu \sim \text{Uniform}(\nu_{\min}, \nu_{\max}). \end{aligned}$$

Recall that  $N$  denotes the total number of observations, i.e.  $\sum_{i=1}^m n_i$  and  $\mathbf{C}_{12} \equiv \text{diag}(b_i)_{1 \leq i \leq N}$ .

### 2.6.1 Approximate Bayesian inference via mean field variational Bayes

A good starting point is to look at the DAG for model (2.16) as illustrated in Figure 2.3. It shares many similarities with the DAG for model (2.11), with the exception of additional parameters  $\mathbf{b}$  and  $\nu$ . This is where the locality property of MFVB approximation comes to the fore through the factorisation afforded by DAG models. As set out in Section 1.4, the distribution of a node within a DAG depends only on the nodes within its Markov blanket. But what impact does this locality property have on the derivation of optimal  $q$ -densities for model (2.16)? Inspecting Figure 2.3 we can see similarities between the DAGs for the Gaussian and Student- $t$  response models. The parameters  $\boldsymbol{\Sigma}^R, \mathbf{a}^R, \sigma_{u\ell}^2$  and  $a_{u\ell}$  play the same role as in the Gaussian response model and therefore we expect to see their optimal  $q$ -density remain unchanged. However, the optimal  $q$ -densities of  $(\boldsymbol{\beta}, \mathbf{u})$  and  $\sigma_\varepsilon^2$

have dependencies on the parameter  $\mathbf{b}$  and therefore we expect to see discrepancies of these optimal  $q$ -densities between the Gaussian and Student- $t$  response models. Ultimately, the locality property of MFVB not only affords considerable simplification for the model at hand, but also allows MFVB results for one model to be transferred to another.

Standard manipulations lead to the following full conditionals for parameters  $\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2, b_i$  and  $\nu$ :

$$\begin{aligned} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} &\sim N \left( (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y}, (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \right), \\ \sigma_\varepsilon^2 |_{\text{rest}} &\sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), \frac{1}{2} (\mathbf{y} - \mathbf{C}\mathbf{v})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C}\mathbf{v}) + a_\varepsilon^{-1} \right), \\ b_i |_{\text{rest}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2} (\nu + 1), \frac{1}{2} \sigma_\varepsilon^{-2} \left\{ (\mathbf{y} - \mathbf{C}\mathbf{v})_i^\top (\mathbf{y} - \mathbf{C}\mathbf{v})_i + \nu \right\} \right), \\ p(\nu |_{\text{rest}}) &= \frac{1}{(\nu_{\max} - \nu_{\min})} \prod_{i=1}^N \left\{ \frac{(0.5\nu)^{0.5\nu}}{\Gamma(0.5\nu)} (b_i)^{-\frac{1}{2}(\nu+2)} \exp \left( -\frac{\nu}{2b_i} \right) \right\}. \end{aligned}$$

Note that the full conditional of  $\nu$  is of non-standard form.

Similar to Subsection 2.5.1, we impose the following product restriction on the approximating  $q$ -densities:

$$\begin{aligned} q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \boldsymbol{\Sigma}^R, \sigma_u^2, \nu, \mathbf{b}) &= q(\boldsymbol{\beta}, \mathbf{u}) q(a_\varepsilon) q(\boldsymbol{\Sigma}^R) q(\sigma_\varepsilon^2) q(\nu) \quad (2.17) \\ &\times \left\{ \prod_{r=1}^{q^R} q(a_r^R) \right\} \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\} \left\{ \sum_{i=1}^N q(b_i) \right\}. \end{aligned}$$

As shown through the derivations in Appendix 2.B, the optimal  $q$ -densities admit the following forms:

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}) &\text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function,} \\ q^*(\sigma_\varepsilon^2) &\text{ is the Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)} \right) \text{ density function,} \\ q^*(a_\varepsilon) &\text{ is the Inverse-Gamma}(1, B_{q(a_\varepsilon)}) \text{ density function,} \\ q^*(\sigma_{u\ell}^2) &\text{ is the Inverse-Gamma} \left( \frac{1}{2} (q_\ell^G + 1), B_{q(\sigma_{u\ell}^2)} \right) \text{ density function,} \\ q^*(a_{u\ell}) &\text{ is the Inverse-Gamma}(1, B_{q(a_{u\ell})}) \text{ density function,} \\ q^*(\boldsymbol{\Sigma}^R) &\text{ is the Inverse-Wishart} \left( \nu + m + q^R + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} \right) \text{ density function,} \\ q^*(a_r^R) &\text{ is the Inverse-Gamma} \left( \frac{1}{2} (\nu + q^R), B_{q(a_r^R)} \right) \text{ density function,} \\ q^*(b_i) &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2} (\mu_{q(\nu)} + 1), B_{q(b_i)} \right) \text{ density function and} \\ q^*(\nu) &= \frac{\exp \left[ N \left\{ \frac{1}{2} \nu \log(\nu/2) - \log \Gamma(\nu/2) \right\} - \frac{1}{2} \nu C_1 \right]}{\mathcal{F}(0, N, C_1, \nu_{\min}, \nu_{\max})}, \quad \nu_{\min} \leq \nu \leq \nu_{\max}, \end{aligned}$$

for parameters  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , the mean vector and covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,

$B_{q(\sigma_\varepsilon^2)}$ , the scale parameter of  $q^*(\sigma_\varepsilon^2)$ ,  $B_{q(a_\varepsilon)}$ , the scale parameter of  $q^*(a_\varepsilon)$ ,  $B_{q(\sigma_{u\ell}^2)}$ , the scale parameter of  $q^*(\sigma_{u\ell}^2)$ ,  $fB_{q(a_{u\ell})}$ , the scale parameter of  $q^*(a_{u\ell})$ ,  $B_{q(a_r^R)}$ , the scale parameter of  $q^*(a_r^R)$  and  $\mathbf{B}_{q(\Sigma^R)}$ , the scale parameter of  $q^*(\Sigma^R)$ .  $B_{q(b_i)}$ , the scale parameter of  $q^*(b_i)$ . The last density uses the definition:  $C_1 \equiv \sum_{i=1}^N \{\mu_{q(\log(b_i))} + \mu_{q(1/b_i)}\}$ .

The MFVB approximation involves standard density functions except for the parameter  $\nu$ . The optimal  $q$ -density  $q^*(\nu)$  involves a non-analytic integral (Wand *et al.*, 2012)

$$\mathcal{F}(p, q, r, s, t) \equiv \int_s^t x^p \exp \left[ q \left\{ \frac{1}{2} x \log(x/2) - \log \Gamma(x/2) \right\} - \frac{1}{2} r x \right] dx,$$

where  $p \geq 0$  and  $q, r, s, t > 0$ , and requires univariate quadrature to solve numerically. A simple and effective form of univariate quadrature is the trapezoidal rule, which can be made arbitrarily accurate as the number of trapezoidal elements increases.

Algorithm 3 gives details of an iterative scheme for finding optimal  $q$ -densities for all key model parameters under product restriction (2.17). Convergence of Algorithm 3 can be monitored using successive values of the lower bound on the marginal log-likelihood  $\log \underline{p}(\mathbf{y}; q)$  and is given as:

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) = & \frac{1}{2} q^R (\nu + q^R - 1) \log(2\nu) - \frac{1}{2} \log(2\pi) \sum_{i=1}^m n_i - \left( \frac{1}{2} q^R + L + 1 \right) \log(\pi) \\ & - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) \\ & + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| - \log(\mathcal{C}_{q^R, \nu + q^R - 1}) + \log(\mathcal{C}_{q^R, \nu + m + q^R - 1}) \\ & - \frac{1}{2} (\nu + m + q^R - 1) \log |\mathbf{B}_{q(\Sigma^R)}| \\ & + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^G + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\ & + \log \Gamma \left\{ \frac{1}{2} (\sum_{i=1}^m n_i + 1) \right\} - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_\varepsilon^2)}) \\ & - \sum_{r=1}^{q^R} \log(A_{Rr}) + q^R \log \Gamma \left\{ \frac{1}{2} (\nu + q^R) \right\} \\ & - \sum_{\ell=1}^L \log(A_{u\ell}) - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\ & + \sum_{r=1}^{q^R} \nu (\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} \mu_{q(1/a_r^R)} - \frac{1}{2} (\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) \\ & - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)} \\ & - \frac{1}{2} \mu_{q(\nu+1)} \sum_{i=1}^N \log B_{q(b_i)} + \log \Gamma \left\{ \frac{1}{2} (\mu_{q(\nu)} + 1) \right\} + \sum_{i=1}^N B_{q(b_i)} \mu_{q(1/b_i)} \\ & - \log(\nu_{\max} - \nu_{\min}) - \frac{1}{2} \mu_{q(\nu)} \sum_{i=1}^N \mu_{q(\log(b_i))} + \log \mathcal{F}(0, \sum_{i=1}^m n_i, C_1, \nu_{\max}, \nu_{\min}), \end{aligned}$$

where  $\boldsymbol{\mu}_{q(\beta)}$  is defined to be the subvector of  $\boldsymbol{\mu}_{q(\beta, \mathbf{u})}$  corresponding to  $\beta$  and similarly,  $\boldsymbol{\Sigma}_{q(\beta)}$  is the sub-matrix of  $\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}$  corresponding to  $\beta$ . The terms  $\mathcal{C}_{q^R, \nu + q^R - 1}$  and  $\mathcal{C}_{q^R, \nu + m + q^R - 1}$  are defined in Table 1.1. Details on the derivation for the optimal  $q$ -densities and lower bound expression are deferred to Appendix 2.B.

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mu_{q(1/\sigma_{u_\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\mu_{q(1/b_i)} > 0$ ,  $1 \leq i \leq N$  and  $\mu_{q(\nu)} > 0$ .

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega^{-1})} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} \left( \mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\Sigma^R)^{-1})} \end{bmatrix}$$

$$N \leftarrow \sum_{i=1}^m n_i \quad ; \quad \text{Gaussian: } \Psi = \mathbf{I} \quad ; \quad \text{Student-}t: \Psi \leftarrow \text{diag}_{1 \leq i \leq N} (\mu_{q(1/b_i)})$$

**1. Update multivariate normal  $q^*(\beta, \mathbf{u})$  parameters:**

$$\Sigma_{q(\beta, \mathbf{u})} \leftarrow (\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \Psi \mathbf{C} + \mathbf{M}_{q(\Omega)^{-1}})^{-1}$$

$$\mu_{q(\beta, \mathbf{u})} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top \Psi \mathbf{y}$$

**2. Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  and inverse-Gamma  $q^*(a_\varepsilon)$  parameters:**

$$B_{q(\sigma_\varepsilon^2)} \leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left\{ (\mathbf{y} - \mathbf{C} \mu_{q(\beta, \mathbf{u})})^\top \Psi (\mathbf{y} - \mathbf{C} \mu_{q(\beta, \mathbf{u})}) + \text{tr}(\mathbf{C}^\top \Psi \mathbf{C} \Sigma_{q(\beta, \mathbf{u})}) \right\}$$

$$\mu_{q(1/\sigma_\varepsilon^2)} \leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2})$$

**3. For Student- $t$  response model only, update inverse-Gamma  $q^*(b_i)$  and  $q^*(\nu)$  parameters:**

For  $1 \leq i \leq N$ :

$$B_{q(b_i)} \leftarrow \frac{1}{2} \left[ \mu_{q(\nu)} + \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C} \mu_{q(\beta, \mathbf{u})})_i^2 + (\mathbf{C} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top)_{ii} \right\} \right]$$

$$\mu_{q(1/b_i)} \leftarrow \frac{1}{2} (\mu_{q(\nu)} + 1) / B_{q(b_i)}$$

$$\mu_{q(\log b_i)} \leftarrow \log(B_{q(b_i)}) - \text{digamma} \left\{ \frac{1}{2} (\mu_{q(\nu)} + 1) \right\}$$

$$C_1 \leftarrow \sum_{i=1}^N \left\{ \mu_{q(\log(b_i))} + \mu_{q(1/b_i)} \right\}$$

$$\mu_{q(\nu)} \leftarrow \exp \left\{ \log \mathcal{F}(1, N, C_1, \nu_{\min}, \nu_{\max}) - \log \mathcal{F}(0, N, C_1, \nu_{\min}, \nu_{\max}) \right\}$$

**4. Update inverse-Gamma  $q^*(a_{u_\ell})$  and inverse-Gamma  $q^*(\sigma_{u_\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\mu_{q(1/a_{u_\ell})} \leftarrow 1 / (\mu_{q(1/\sigma_{u_\ell}^2)} + A_{u_\ell}^{-2})$$

$$\mu_{q(\sigma_{u_\ell}^2)} \leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}_\ell^G)}) + 2 \mu_{q(1/a_{u_\ell})}}$$

**5. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\Sigma^R)$  parameters:**

 For  $r = 1, \dots, q^R$ :

$$B_{q(a_r^R)} \leftarrow \nu(M_{q((\Sigma^R)^{-1})})_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a_r^R)} \leftarrow \frac{1}{2}(\nu + q^R)/B_{q(a_r^R)}$$

$$B_{q(\Sigma^R)} \leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \Sigma_{q(\mathbf{u}_i^R)}) + 2\nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)})$$

$$M_{q((\Sigma^R)^{-1})} \leftarrow (\nu + m + q^R - 1) B_{q(\Sigma^R)}^{-1}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**


---

Algorithm 3: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed models with the Gaussian and Student- $t$  responses.

## 2.7 Bernoulli semiparametric mixed models

We now consider the case where the response distribution of  $\mathbf{y}$  is Bernoulli distributed, the logistic mixed model

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u} &\sim \text{Bernoulli}(\text{logit}^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})), \\ \text{where} \quad \text{logit}^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) &= \frac{e^{\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}}}{\mathbf{1} + e^{\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}}} \end{aligned} \quad (2.18)$$

is appropriate. The notation  $v \sim \text{Bernoulli}(p)$  used in (2.18) is shorthand for the entries of the random vector  $v$  having independent Bernoulli distributions with parameters corresponding to the entries of  $p$ .

The notational infrastructure described in Section 2.5 for the Gaussian response model can be directly transferred to the Bernoulli response case. The only change is the response distribution specification and omission of parameters  $\sigma_\varepsilon^2$  and  $a_\varepsilon$ . The two-level Bayesian Bernoulli semiparametric mixed model stated in full is:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim \text{Bernoulli}(\text{logit}^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), \quad (2.19) \\ \mathbf{u} | \sigma_{u_1}^2, \dots, \sigma_{u_L}^2, \Sigma^R &\sim N\left(\mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \Sigma^R \end{bmatrix}\right), \\ \Sigma^R | a_1^R, \dots, a_{q^R}^R &\sim \text{Inverse-Wishart}\left(\nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R)\right), \\ a_r^R &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{Rr}^2\right), \quad r = 1, \dots, q^R, \\ \sigma_{u\ell}^2 | a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{u\ell}\right), \quad a_{u\ell} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{u\ell}^2\right). \end{aligned}$$

### 2.7.1 Approximate Bayesian inference via mean field variational Bayes

We once again utilise the locality property of MFVB approximation to derive the necessary optimal  $q$ -densities for model (2.19). Inspection of Figure 2.4 reviews great similarities between the DAGs for the Gaussian and Bernoulli response models. The only change that is required would be for  $(\boldsymbol{\beta}, \mathbf{u})$  since the response vector  $\mathbf{y}$  is now Bernoulli distributed and hence the optimal  $q$ -density  $q^*(\boldsymbol{\beta}, \mathbf{u})$  is the focus of this subsection.

The marginal likelihood is

$$p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}) = \exp \left[ \mathbf{y}^\top (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top \log \{ \mathbf{1} + \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \} \right].$$

The joint density function of  $(\boldsymbol{\beta}, \mathbf{u})$  satisfies

$$\begin{aligned} p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}) &\propto \exp \left[ \mathbf{y}^\top (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top \log \{ \mathbf{1} + \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \} \right] \\ &\quad \times \prod_{i=1}^m \exp \left[ -\frac{1}{2} (\mathbf{u}^R)^\top \{ \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^R)^{-1} \} \mathbf{u}^R \right] \\ &\quad \times \prod_{\ell=1}^L \exp \left\{ -\frac{1}{2} (\mathbf{u}_\ell^G)^\top (\sigma_{u_\ell}^2 \mathbf{I}_{q_\ell})^{-1} \mathbf{u}_\ell^G \right\} \exp \left\{ -\frac{1}{2} \boldsymbol{\beta}^\top (\sigma_\beta^2 \mathbf{I}_p)^{-1} \boldsymbol{\beta} \right\}. \end{aligned}$$

This is a non-standard form and poses tractability problems with regard to approximate inference for  $(\boldsymbol{\beta}, \mathbf{u})$ . The term  $-\log(1 + e^x)$  is a convex function of  $x$  and we propose to transform this convex term to a simple quadratic function  $f(x; \boldsymbol{\xi})$ , a trick first introduced in Jaakkola and Jordan (1997). Different values of  $\boldsymbol{\xi}$  correspond to different parabolas, all of which is smaller than  $-\log(1 + e^x)$ , and thus we have the following representation of  $-\log(1 + e^x)$  as the maxima of a family of parabolas (Ormerod and Wand, 2010):

$$-\log(1 + e^x) = \max_{\xi \in \mathbb{R}} \left\{ \lambda(\xi) x^2 - \frac{1}{2} x + \psi(\xi) \right\} \text{ for all } x \in \mathbb{R}, \quad (2.20)$$

$$\text{where } \lambda(\xi) \equiv -\tanh(\xi/2)/(4\xi)$$

$$\text{and } \psi(\xi) \equiv \xi/2 - \log(1 + e^\xi) + \xi \tanh(\xi/2)/4.$$

Figure 2.6 is the graphical representation of (2.20), in which the function  $-\log(1 + e^x)$  is seen to be the maximum of a family of parabolas. It follows from (2.20) that

$$\begin{aligned} \mathbf{1}^\top \log \{ \mathbf{1} + \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \} &\geq \mathbf{1}^\top \left\{ \lambda(\boldsymbol{\xi}) \odot (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})^2 - \frac{1}{2} (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) + \psi(\boldsymbol{\xi}) \right\}, \\ &= (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \\ &\quad - \frac{1}{2} \mathbf{1}^\top (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) + \mathbf{1}^\top \psi(\boldsymbol{\xi}), \end{aligned}$$

where  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_N)$  is an  $(\sum_{i=1}^m n_i) \times 1$  vector of variational parameters. Recall that



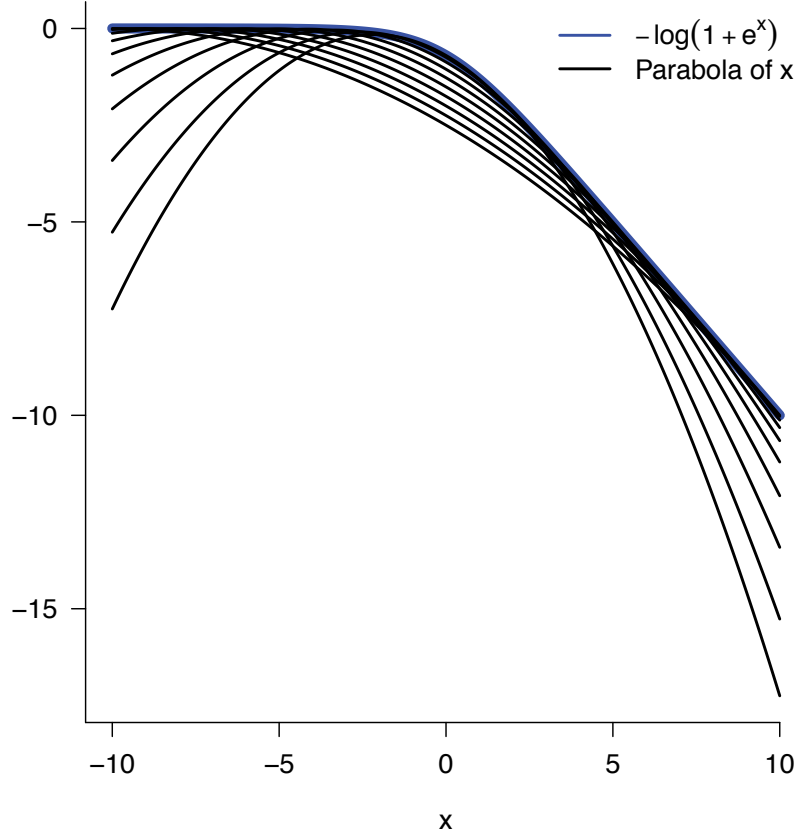


Figure 2.6: Variational representation of the function  $-\log(1 + e^x)$  as the maximum of a family of parabolas, corresponding to (2.20).

$\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$  and  $\mathbf{v} = [\boldsymbol{\beta}^\top \ \mathbf{u}^\top]^\top$ , the lower bound on  $p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$  is

$$\underline{p}(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) \propto \exp(\mathbf{y}^\top \mathbf{C} \mathbf{v} - [\mathbf{v}^\top \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} \mathbf{v} - \frac{1}{2} \mathbf{1}^\top \mathbf{C} \mathbf{v} + \mathbf{1}^\top \psi(\boldsymbol{\xi})] - \frac{1}{2} \mathbf{v}^\top \boldsymbol{\Omega}^{-1} \mathbf{v}),$$

where  $\boldsymbol{\Omega}$  is as defined in Section 2.4. Noting that since the right hand side of the above equation is a quadratic form and, by completing the square in the usual way to identify the mean and covariance, we approximate the full conditional of  $(\boldsymbol{\beta}, \mathbf{u})$  by a multivariate normal distribution

$$\boldsymbol{\beta}, \mathbf{u} | \mathbf{y}; \boldsymbol{\xi} \sim N\left( (2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \mathbf{C}^\top (\mathbf{y} - \frac{1}{2} \mathbf{1}), (2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \right).$$

The Bernoulli response model entails a variant of variational approximation, namely the *tangent transform approach* of Jaakkola and Jordan (1997). As described above, we succeed in simplifying the convex function of  $x$  with a simple quadratic function. The only trade-off is to include the additional variational parameter vector  $\boldsymbol{\xi}$  and to obtain

the lower bound with respect to  $\boldsymbol{\xi}$ . The Jaakkola and Jordan (1997) trick is one of a few approaches that has been proposed in the variational approximation literature for handling binary response regression models. Others include the use of probit auxiliary variable trick of Albert and Chib (1993) and Gaussian variational approximation (e.g. Knowles and Minka, 2011; Ormerod and Wand, 2012). The Gaussian variational approximation approach emerges as a fast, deterministic alternative to Laplace approximation for fitting of grouped data generalised linear mixed models. The approach is conceptually simple and its derivation requires application of Jensen's inequality to the log-likelihood to obtain a variational lower bound. Various trade-offs are involved with the choice amongst these options. For example, Gaussian variational approximation requires univariate numerical integration whereas Algorithm 4 has only purely algebraic updates.

Similar to Subsection 2.5.1, the  $q$ -densities are subject to the following product form:

$$q(\boldsymbol{\beta}, \mathbf{u}) q(\boldsymbol{\Sigma}^R) \left\{ \prod_{r=1}^{q^R} q(a_r^R) \right\} \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\}. \quad (2.21)$$

As shown through the derivations in Appendix 2.C, the optimal  $q$ -densities admit the following forms:

$$\begin{aligned} \boldsymbol{\xi} &\leftarrow \sqrt{\text{diagonal} \left\{ \mathbf{C} \left( \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^\top \right) \right\} \mathbf{C}^\top}, \\ q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) &\text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}) \text{ density function,} \\ q^*(\sigma_{u\ell}^2) &\text{ is the Inverse-Gamma} \left( \frac{1}{2}(q_\ell^G + 1), B_{q(\sigma_{u\ell}^2)} \right) \text{ density function,} \\ q^*(a_{u\ell}) &\text{ is the Inverse-Gamma}(1, B_{q(a_{u\ell})}) \text{ density function,} \\ q^*(\boldsymbol{\Sigma}^R) &\text{ is the Inverse-Wishart} \left( \nu + m + q^R + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} \right) \text{ density function,} \\ q^*(a_r^R) &\text{ is the Inverse-Gamma} \left( \frac{1}{2}(\nu + q^R), B_{q(a_r^R)} \right) \text{ density function,} \end{aligned}$$

for parameters  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , the mean vector and covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $B_{q(\sigma_{u\ell}^2)}$ , the scale parameter of  $q^*(\sigma_{u\ell}^2)$ ,  $B_{q(a_{u\ell})}$ , the scale parameter of  $q^*(a_{u\ell})$ ,  $B_{q(a_r^R)}$ , the scale parameter of  $q^*(a_r^R)$  and  $\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}$ , the scale parameter of  $q^*(\boldsymbol{\Sigma}^R)$ .

Algorithm 4 gives details of an iterative scheme for finding optimal  $q$ -densities for all key model parameters under product restriction (2.21). Convergence of Algorithm 4 can be monitored using successive values of the lower bound on the marginal log-likelihood and is given as:

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^R (\nu + q^R - 1) \log(2\nu) - \left( \frac{1}{2} q^R + L \right) \log(\pi) - \lambda(\boldsymbol{\xi})^\top (\boldsymbol{\xi}^2) \\ &\quad + \mathbf{1}^\top \psi(\boldsymbol{\xi}) + (\mathbf{y} - \frac{1}{2} \mathbf{1})^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}| \\ &\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) \end{aligned}$$

$$\begin{aligned}
 & -\log(\mathcal{C}_{q^R, \nu+q^R-1}) + \log(\mathcal{C}_{q^R, \nu+m+q^R-1}) - \frac{1}{2}(\nu + m + q^R - 1) \log|B_{q(\Sigma^R)}| \\
 & + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2}(q_\ell^G + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\
 & - \sum_{r=1}^{q^R} \log(A_{Rr}) + q^R \log \Gamma \left\{ \frac{1}{2}(\nu + q^R) \right\} \\
 & + \sum_{r=1}^{q^R} \nu (\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} \mu_{q(1/a_r^R)} \\
 & - \frac{1}{2}(\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) - \sum_{\ell=1}^L \log(A_{u\ell}) \\
 & - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)},
 \end{aligned}$$

where  $\boldsymbol{\mu}_{q(\boldsymbol{\beta})}$  is defined to be the subvector of  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$  and similarly,  $\Sigma_{q(\boldsymbol{\beta})}$  is the sub-matrix of  $\Sigma_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$ . The terms  $\mathcal{C}_{q^R, \nu+q^R-1}$  and  $\mathcal{C}_{q^R, \nu+m+q^R-1}$  are defined in Table 1.1. Details on the derivation for the optimal  $q$ -densities and lower bound expression are deferred to Appendix 2.C.

## 2.8 Poisson semiparametric mixed models

Finally, we consider the distribution of response vector  $\mathbf{y}$  as Poisson, the full two-level Bayesian Poisson semiparametric mixed model is

$$\begin{aligned}
 \mathbf{y} | \boldsymbol{\beta}, \mathbf{u} & \sim \text{Poisson}(\exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), \quad (2.22) \\
 \mathbf{u} | \sigma_{u1}^2, \dots, \sigma_{uL}^2, \Sigma^R & \sim N\left(\mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \Sigma^R \end{bmatrix}\right), \\
 \Sigma^R | a_1^R, \dots, a_{q^R}^R & \sim \text{Inverse-Wishart}\left(\nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R)\right), \\
 a_r^R & \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{Rr}^2\right), \quad r = 1, \dots, q^R, \\
 \sigma_{u\ell}^2 | a_{u\ell} & \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{u\ell}\right), \quad a_{u\ell} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{u\ell}^2\right).
 \end{aligned}$$

Figure 2.4 provides a DAG representation of (2.19) and (2.22). Observed data are indicated by the shaded node and random parameters are represented by the open nodes. The visual representation shows that the Bernoulli and Poisson response models have all parts of their graphs in common. The locality property of MFVB approximation suggests that the algorithms for the two models share some of the updates.

### 2.8.1 Approximate Bayesian inference via non-conjugate variational message passing

In generalised response regression, the Poisson distribution is often bracketed with the Bernoulli distribution since both belong to the one-parameter exponential family. However, variational approximations for the Poisson response models are not the same as those with the Bernoulli response models. For the Bayesian logistic mixed models, Jaakkola and

**Initialise:**  $\nu = 2$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$  and  $\boldsymbol{\xi}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega^{-1})} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} \left( \mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\Sigma^R)^{-1})} \end{bmatrix}$$

**1. Update multivariate normal  $q^*(\boldsymbol{\beta}, \mathbf{u})$  and  $\boldsymbol{\xi}$  parameters:**

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \left[ 2 \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \mathbf{M}_{q(\Omega^{-1})} \right]^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top (\mathbf{y} - \tfrac{1}{2} \mathbf{1})$$

$$\boldsymbol{\xi} \leftarrow \sqrt{\text{diagonal} \left\{ \mathbf{C} \left( \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^\top \right) \mathbf{C}^\top \right\}}$$

**2. Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\begin{aligned} \mu_{q(1/a_{u\ell})} &\leftarrow 1 / (\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2}) \\ \mu_{q(\sigma_{u\ell}^2)} &\leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)}) + 2 \mu_{q(1/a_{u\ell})}} \end{aligned}$$

**3. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\Sigma^R)$  parameters:**

For  $r = 1, \dots, q^R$ :

$$B_{q(a_r^R)} \leftarrow \nu (\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a_r^R)} \leftarrow \frac{1}{2} (\nu + q^R) / B_{q(a_r^R)}$$

$$\mathbf{B}_{q(\Sigma^R)} \leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}) + 2 \nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)})$$

$$\mathbf{M}_{q((\Sigma^R)^{-1})} \leftarrow (\nu + m + q^R - 1) \mathbf{B}_{q(\Sigma^R)}^{-1}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**

---

Algorithm 4: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Bernoulli response.

Jordan (1997) derive a lower bound on the marginal likelihood that leads to tractable approximate inference. The Albert and Chib (1993) auxiliary variable representation of Bayesian probit regression leads to a different type of variational approximation for the binary outcomes. There do not appear to be analogues of these approaches for the Poisson regression and therefore a different technique is explored.

We begin with the MFVB approximation such that the actual joint posterior density function is approximated by the product restricted form

$$p(\boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}^R, \mathbf{a}^R, \boldsymbol{\sigma}_u^2, \mathbf{a}_u | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) q(\mathbf{a}^R, \mathbf{a}_u). \quad (2.23)$$

However, the optimal  $q$ -density function  $q^*(\boldsymbol{\beta}, \mathbf{u})$  under (2.23) satisfies

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \propto \exp [E_{q(-(\boldsymbol{\beta}, \mathbf{u}))} \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \text{Markov blanket of } (\boldsymbol{\beta}, \mathbf{u})) \}],$$

which involves non-closed form multivariate integrals and therefore is not tractable. A remedy in this instance is to postulate  $q(\boldsymbol{\beta}, \mathbf{u})$  to be a special case of the exponential family density function, in this case multivariate normal

$$p(\boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}^R, \mathbf{a}^R, \boldsymbol{\sigma}_u^2, \mathbf{a}_u | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) q(\boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) q(\mathbf{a}^R, \mathbf{a}_u),$$

where

$$q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function.} \quad (2.24)$$

Recall that  $\mathbf{v} = \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}$  and express (2.24) as an exponential family form gives

$$p(\mathbf{v}) = \exp \left\{ \mathbf{T}(\mathbf{v})^\top \boldsymbol{\eta} - A(\boldsymbol{\eta}) - \frac{d}{2} \log(2\pi) \right\},$$

where

$$\begin{aligned} \mathbf{T}(\mathbf{v}) &\equiv \begin{bmatrix} \mathbf{v} \\ \text{vech}(\mathbf{v}\mathbf{v}^\top) \end{bmatrix}, & \boldsymbol{\eta} &\equiv \begin{bmatrix} \boldsymbol{\Sigma}^{-1} \mathbf{u} \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} \\ \text{and } A(\boldsymbol{\eta}) &= -\frac{1}{4} \boldsymbol{\eta}_1^\top \left\{ \text{vec}^{-1}(\mathbf{D}_2^{+\top} \boldsymbol{\eta}_2) \right\}^{-1} - \frac{1}{2} \log \left| -2 \text{vec}^{-1}(\mathbf{D}_d^{+\top} \boldsymbol{\eta}_2) \right| \end{aligned}$$

are the natural statistic, natural parameter and log-partition function respectively. The optimal  $q$ -density function  $q^*(\boldsymbol{\beta}, \mathbf{u})$  can instead be found using

$$\boldsymbol{\eta} \leftarrow [\text{var} \{ \mathbf{T}(\mathbf{v}) \}]^{-1} \{ \mathbf{D}_\eta E_{q(\boldsymbol{\Omega}, \mathbf{v})} \{ \log p(\mathbf{y}, \boldsymbol{\Omega}, \mathbf{v}) \} \}, \quad (2.25)$$

where  $\boldsymbol{\Omega} \in \{ \boldsymbol{\Sigma}^R, \mathbf{a}^R, \boldsymbol{\sigma}_u^2, \mathbf{a}_u \}$ . This approach, first introduced in Knowles and Minka

(2011), has been labelled as the *non-conjugate variational message passing* since it offers a way of circumventing non-conjugacies in ordinary MFVB. The exponential family distribution parameters are chosen via fixed-point iteration, which we describe in Theorem 2.8.1 following a straight adaption from Rhode and Wand (2015). Wand (2014) provides fully simplified fixed point update expressions for the multivariate normal  $q$ -density parameters, in terms of derivatives with respect to the corresponding mean and covariance matrix:

$$\begin{aligned}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} (\mathbf{D}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \mathbb{E}_{q(\boldsymbol{\Omega}, \boldsymbol{\beta}, \mathbf{u})} \{\log p(\mathbf{y}, \boldsymbol{\Omega}, \boldsymbol{\beta}, \mathbf{u})\})^\top, \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow - \left( \mathbf{H}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \mathbb{E}_{q(\boldsymbol{\Omega}, \boldsymbol{\beta}, \mathbf{u})} \{\log p(\mathbf{y}, \boldsymbol{\Omega}, \boldsymbol{\beta}, \mathbf{u})\} \right)^{-1},\end{aligned}\quad (2.26)$$

where  $\mathbf{H}$  denotes the Hessian matrix as defined in Definition 1.9.2. The term  $\mathbf{D}_{\mathbf{x}}f$  is the derivative factor of  $f$  with respect to  $\mathbf{x}$ .

**Theorem 2.8.1.** The aim of fixed-point iterations is to find points, denoted by  $\mathbf{x}^*$ , that satisfy the following stationary point condition

$$\mathbf{D}_{\mathbf{x}}f(\mathbf{x})^\top = \mathbf{0}, \quad (2.27)$$

assuming that  $\mathbf{D}_{\mathbf{x}}f(\mathbf{x})$  exists. Such points are considered as the maxima or minima of  $f$ . Equation (2.27) can be rewritten in the form  $\mathbf{x} = \mathbf{g}(\mathbf{x})$  for some function  $\mathbf{g} : D \subseteq \mathbb{R}^d \leftarrow \mathbb{R}^d$ . Given  $\mathbf{g}$ , fixed-point iterations simply involve repeated evaluation of  $\mathbf{g}$  as given in Algorithm 5.

---

**Initialise:**  $\mathbf{x} \leftarrow \mathbf{x}_{\text{init}}$  for some  $\mathbf{x}_{\text{init}} \in D$ .

**Cycle:**

$$\mathbf{x} \leftarrow \mathbf{g}(\mathbf{x})$$

**until convergence.**

---

Algorithm 5: The fixed-point iteration algorithm in generic form.

Similar to Subsection 2.5.1, the  $q$ -densities are subject to the following product form:

$$q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^{\mathbf{R}}, \mathbf{a}_u, \boldsymbol{\Sigma}^{\mathbf{R}}, \sigma_u^2) = q(\boldsymbol{\beta}, \mathbf{u}) q(\boldsymbol{\Sigma}^{\mathbf{R}}) \left\{ \prod_{r=1}^{q^{\mathbf{R}}} q(a_r^{\mathbf{R}}) \right\} \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\} \quad (2.28)$$

As shown through the derivations in Appendix 2.D, the optimal  $q$ -densities admit the

following forms:

$$\begin{aligned}
 q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) & \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function,} \\
 q^*(\sigma_{u\ell}^2) & \text{ is the Inverse-Gamma } \left( \frac{1}{2}(q_\ell^G + 1), B_{q(\sigma_{u\ell}^2)} \right) \text{ density function,} \\
 q^*(a_{u\ell}) & \text{ is the Inverse-Gamma } (1, B_{q(a_{u\ell})}) \text{ density function,} \\
 q^*(\boldsymbol{\Sigma}^R) & \text{ is the Inverse-Wishart } \left( \nu + m + q^R + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} \right) \text{ density function,} \\
 q^*(a_r^R) & \text{ is the Inverse-Gamma } \left( \frac{1}{2}(\nu + q^R), B_{q(a_r^R)} \right) \text{ density function,}
 \end{aligned}$$

for parameters  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , the mean vector and covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $B_{q(\sigma_{u\ell}^2)}$ , the scale parameter of  $q^*(\sigma_{u\ell}^2)$ ,  $B_{q(a_{u\ell})}$ , the scale parameter of  $q^*(a_{u\ell})$ ,  $B_{q(a_r^R)}$ , the scale parameter of  $q^*(a_r^R)$  and  $\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}$ , the scale parameter of  $q^*(\boldsymbol{\Sigma}^R)$ .

The interdependencies between the parameters in these optimal  $q$ -density functions, combined with the updates for  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  given by (2.26), give rise to an iterative scheme for their solution and is encompassed in Algorithm 6. Algorithm 6 also uses the variational lower bound on the marginal log-likelihood. For model (2.22) under product restriction (2.28) it has the explicit expression as follows:

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^R (\nu + q^R - 1) \log(2\nu) - \left( \frac{1}{2} q^R + L \right) \log(\pi) \\
 &+ \mathbf{y}^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} - \mathbf{1}^\top \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \mathbf{C}^\top \right\} \\
 &- \mathbf{1}^\top \log(\mathbf{y}!) + \left( \mathbf{y} - \frac{1}{2} \mathbf{1} \right)^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}| \\
 &- \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right) \\
 &- \log(\mathcal{C}_{q^R, \nu + q^R - 1}) + \log(\mathcal{C}_{q^R, \nu + m + q^R - 1}) \\
 &- \frac{1}{2} (\nu + m + q^R - 1) \log |B_{q(\boldsymbol{\Sigma}^R)}| \\
 &+ \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^G + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\
 &+ q^R \log \Gamma \left\{ \frac{1}{2} (\nu + q^R) \right\} + \sum_{r=1}^{q^R} \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})})_{rr} \mu_{q(1/a_r^R)} \\
 &- \frac{1}{2} (\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) - \sum_{\ell=1}^L \log(A_{u\ell}) \\
 &- \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\
 &+ \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) - \sum_{r=1}^{q^R} \log(A_{Rr}),
 \end{aligned}$$

where  $\boldsymbol{\mu}_{q(\boldsymbol{\beta})}$  is defined to be the subvector of  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$  and similarly,  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}$  is the sub-matrix of  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$ . The terms  $\mathcal{C}_{q^R, \nu + q^R - 1}$  and  $\mathcal{C}_{q^R, \nu + m + q^R - 1}$  are defined in Table 1.1. Details on the derivation for the optimal  $q$ -densities and lower bound expression are deferred to Appendix 2.D. It is worth mentioning that, unlike ordinary MFVB, there is no guarantee that there will be an increase at each iteration for this variational lower bound. Therefore, we continue the cycle of updates once the absolute relative change in its logarithm falls below a negligible amount.

**Initialise:**  $\nu = 2$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\boldsymbol{\mu}_{q(\beta, \mathbf{u})}$ , a  $(p + q^R m + \sum_{\ell=1}^L q_\ell^G) \times 1$  vector and  $\mathbf{w}_{q(\beta, \mathbf{u})}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega^{-1})} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} \left( \mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\Sigma^R)^{-1})} \end{bmatrix}$$

**Update multivariate normal  $q^*(\beta, \mathbf{u})$  and  $\mathbf{w}_{q(\beta, \mathbf{u})}$  parameters:**

$$\begin{aligned} \Sigma_{q(\beta, \mathbf{u})} &\leftarrow \left\{ \mathbf{C}^\top \text{diag}(\mathbf{w}_{q(\beta, \mathbf{u})}) \mathbf{C} + \mathbf{M}_{q(\Omega^{-1})} \right\}^{-1} \\ \boldsymbol{\mu}_{q(\beta, \mathbf{u})} &\leftarrow \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \Sigma_{q(\beta, \mathbf{u})} \left\{ \mathbf{C}^\top (\mathbf{y} - \mathbf{w}_{q(\beta, \mathbf{u})}) - \mathbf{M}_{q(\Omega^{-1})} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \right\} \\ \mathbf{w}_{q(\beta, \mathbf{u})} &\leftarrow \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top) \right\} \end{aligned}$$

**2. Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\begin{aligned} \mu_{q(1/a_{u\ell})} &\leftarrow 1 / (\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2}) \\ \mu_{q(\sigma_{u\ell}^2)} &\leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}_\ell^G)}) + 2 \mu_{q(1/a_{u\ell})}} \end{aligned}$$

**3. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\Sigma^R)$  parameters:**

For  $r = 1, \dots, q^R$ :

$$\begin{aligned} B_{q(a_r^R)} &\leftarrow \nu (\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a_r^R)} \leftarrow \frac{1}{2} (\nu + q^R) / B_{q(a_r^R)} \\ \mathbf{B}_{q(\Sigma^R)} &\leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \Sigma_{q(\mathbf{u}_i^R)}) + 2 \nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)}) \\ \mathbf{M}_{q((\Sigma^R)^{-1})} &\leftarrow (\nu + m + q^R - 1) \mathbf{B}_{q(\Sigma^R)}^{-1} \end{aligned}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**

---

Algorithm 6: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Poisson response.



## 2.9 Displaying approximate posterior means of regression function fits

The majority of this chapter describes the process for obtaining approximating posterior distributions for various mixed model parameters. In this section we translate these into meaningful graphical displays consisting of fitted curves and their corresponding variability bands.

Essentially, we plot for each predictor, the profile of the response outcome with each of the other predictors set at their mean value. In practice, we set up a grid over the domain of each predictor. Take the Gaussian response model with  $L = 1$  as an example, at the global level, this would involve plotting the overall mean of the penalised spline regression function

$$\hat{h}(\bar{x}; s) = \hat{\beta}_0 + \hat{\beta}_x \bar{x} + \hat{\beta}_s s + \sum_{k=1}^{q^G} u_k^G z_k(s).$$

Here we simply ignore the random effects terms that represent group-specific departures from that overall mean  $\hat{h}(\bar{x}; s)$ .

Penalised regression splines lend themselves to fairly straightforward pointwise credible bands as we now describe. Suppose that we want a pointwise credible band for function  $\hat{h}(\bar{x}; s)$  over a grid of  $M$   $s$ -values, i.e.  $\mathbf{g} = (g_{s1}, \dots, g_{sM})$ . Define  $\hat{h}_g$  to be the MFVB approximate posterior mean of the penalised spline function over  $\mathbf{g}$  and  $\text{Cov}_q(\cdot)$  be the  $q$ -density covariance matrix corresponding to  $(\boldsymbol{\beta}, \mathbf{u}^G)$ :

$$\hat{h}_g \equiv \begin{bmatrix} \hat{h}_g(\bar{x}, g_{s1}) \\ \vdots \\ \hat{h}_g(\bar{x}, g_{sM}) \end{bmatrix} = \mathbf{C}_g \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \quad \text{and} \quad \text{Cov}_q \left( \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix} \right) = \begin{bmatrix} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} & \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \\ \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}^\top & \boldsymbol{\Sigma}_{q(\mathbf{u}^G)} \end{bmatrix},$$

where  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}$  is defined to be the subvector of  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  corresponding to  $\boldsymbol{\beta}$  and  $\mathbf{u}^G$ , and the remaining terms are defined analogously. Variability bands in  $\hat{h}_g$  can then be obtained from

$$\text{Var}_q(\hat{h}_g) = \mathbf{C}_g \text{Cov}_q \left( \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix} \right) \mathbf{C}_g^\top,$$

where

$$\mathbf{C}_g \equiv \begin{bmatrix} 1 & \bar{x} & g_{s1} & z_1(g_{s1}) & \cdots & z_{q^G}(g_{s1}) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \bar{x} & g_{sM} & z_1(g_{sM}) & \cdots & z_{q^G}(g_{sM}) \end{bmatrix}.$$

The MFVB approximation of  $\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix}$ , i.e.  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}$  as presented previously, has a multivariate normal distribution. Thus, our problem reduced to finding the 0.025 and 0.975 quantiles of a multivariate normal at each point on our grid. It follows that an approximate  $100(1 - \alpha)\%$  credible interval for  $\hat{h}_g$  is

$$\hat{h}_g \pm z(1 - \alpha/2)\sqrt{\text{Var}_q(\hat{h}_g)},$$

and the display of this fit is achieved by plotting  $\hat{h}_g$  against  $\mathbf{g}$ . The variability estimates in the group-specific mean function estimates can be obtained in a similar vein.

## 2.10 Numerical evaluation

We conducted a series of comprehensive simulation studies to assess the performance of Algorithms 3, 4 and 6 in terms of inferential accuracy, credible interval coverage and computational speed. We generated 30 datasets according to the following simulation settings:

### Gaussian response

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R) x_{ij} + f(s_{ij}), \sigma_\varepsilon^2),$$

$$\text{where } f(s) = 1 - 2.6\phi(s; 0.15, 0.1) - (2.3s - 0.07s^2) + 0.5\{1 - \Phi(s; 0.8, 0.07)\},$$

$$x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad [u_{0i}^R \quad u_{1i}^R]^\top \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}^R).$$

$$\text{True values: } \beta_0 = 0.58, \quad \beta_x = 1.89, \quad \sigma_\varepsilon^2 = 0.1 \quad \text{and} \quad \boldsymbol{\Sigma}^R = \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix}.$$

### Student-*t* response

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2, \nu \stackrel{\text{ind.}}{\sim} t(\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R) x_{ij} + f(s_{ij}), \sigma_\varepsilon^2, \nu),$$

$$\text{where } f(s) = 0.1 + \cos(4\pi s), \quad x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1),$$

$$s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad [u_{0i}^R \quad u_{1i}^R]^\top \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}^R).$$

$$\text{True values: } \nu = 2, \quad \beta_0 = -1.58, \quad \beta_x = 0.30, \quad \sigma_\varepsilon^2 = 0.1 \quad \text{and} \quad \boldsymbol{\Sigma}^R = \begin{bmatrix} 1.91 & 0.34 \\ 0.34 & 3.27 \end{bmatrix}.$$

### Bernoulli response

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R \stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\text{logit}^{-1}\{\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R) x_{ij} + f(s_{ij})\}),$$

$$\text{where } f(s) = 1 - 2.6\phi(s; 0.15, 0.1) - (2.3s - 0.07s^2) + 0.5\{1 - \Phi(s; 0.8, 0.07)\},$$

$$x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad [u_{0i}^{\text{R}} \quad u_{1i}^{\text{R}}]^{\text{T}} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}^{\text{R}}).$$

True values:  $\beta_0 = 0.58, \quad \beta_x = 1.89 \quad \text{and} \quad \boldsymbol{\Sigma}^{\text{R}} = \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix}.$

### Poisson response

$$y_{ij} | \beta_0, \beta_x, u_{0i}^{\text{R}}, u_{1i}^{\text{R}} \stackrel{\text{ind.}}{\sim} \text{Poisson}(\exp\{\beta_0 + u_{0i}^{\text{R}} + (\beta_x + u_{1i}^{\text{R}})x_{ij} + f(s_{ij})\}),$$

where  $f(s) = 0.1 + \cos(4\pi s); \quad x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1),$

$$s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1), \quad [u_{0i}^{\text{R}} \quad u_{1i}^{\text{R}}]^{\text{T}} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}^{\text{R}}).$$

True values:  $\beta_0 = -1.58, \quad \beta_x = 0.30, \quad \sigma_\varepsilon^2 = 0.1 \quad \text{and} \quad \boldsymbol{\Sigma}^{\text{R}} = \begin{bmatrix} 1.91 & 0.34 \\ 0.34 & 3.27 \end{bmatrix}.$

The above settings are for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ . The number of groups  $m$  is 50 and the within-group sample size  $n_i$  ranged between 40 and 50. All hyperparameters are set to 10 000 and lastly we use  $q^{\text{G}} = 25$  interior knots when using the O’Sullivan splines with spacing as described in Wand and Ormerod (2008).

#### 2.10.1 Fitting a Bayesian hierarchical model in Stan

We describe in detail the implementation of model (2.11) in Section 2.5 via Stan. Note that the necessary changes for other response models and more complex models are relatively minor. These steps include writing the model in Stan, using Stan to set up the data and starting values, calling Stan and graphing the results.

```
'data
{
  int <lower=1> n;           int <lower=1> m;
  int <lower=1> numKnots;   int <lower=1> idnum[n];
  int <lower=1> ncX;       matrix[n,ncX] X;
  vector[n] y;           matrix[n,numKnots] ZG;
  real x[n];             real <lower=0> sigmaBeta;
  real <lower=1> Aeps;    real <lower=1> Au;
  real <lower=1> Ar;      vector[2] zeroVec;
}
parameters
{
  vector[ncX] beta;       matrix[2,m] uR;
  vector[numKnots] uG;    cov_matrix[2] SigmaR;
  real <lower=0> sigmauG;  real <lower=0> sigmaEps;
  vector <lower=1e-10> [2] aR;
}
transformed parameters
{
```

## 2.10. NUMERICAL EVALUATION

---

```

vector[n] fmean;          real mu[n];
fmean <- X*beta + ZG*uG;
for (i in 1:n)
  mu[i] <- (fmean[i] + uR[1,idnum[i]] + uR[2,idnum[i]]*x[i]);
}
model
{
  matrix[2,2] rateForWish;

  y ~ normal(mu,sigmaEps);

  for (j in 1:m)
    col(uR,j) ~ multi_normal(zeroVec,SigmaR);

  rateForWish[1,2] <- 0 ; rateForWish[2,1] <- 0 ;
  for (r in 1:2)
  {
    aR[r] ~ inv_gamma(0.5,pow(Ar,-2));
    rateForWish[r,r] <- 4/aR[r];
  }
  SigmaR ~ inv_wishart(3,rateForWish);

  uG ~ normal(0,sigmauG);
  beta ~ normal(0,sigmaBeta);
  sigmaEps ~ cauchy(0,Aeps);
  sigmauG ~ cauchy(0,Au);
}'

```

The first paragraph of the above code specifies the data: the number of observations,  $N$ ; the number of groups,  $m$ ; the response vector,  $\mathbf{y}$ ; the number of spline knots,  $q^G$ ; the group identification number, `idnum`; the respective fixed and random effects design matrices,  $\mathbf{X}$ ,  $\mathbf{Z}^R$  and  $\mathbf{Z}^G$ ; and the hyperparameters,  $\sigma_\beta^2$ ,  $A_\varepsilon$ ,  $A_R$  and  $A_u$ . Data are labelled as integer or real and can be vectors if dimensions are specified. Data can also be constrained, for example all hyperparameters must be positive. The code next introduces the unknowns to be estimated in the model fit. These are the fixed effects vector,  $\boldsymbol{\beta}$ ; the random effects vectors,  $\mathbf{u}^R$  and  $\mathbf{u}^G$ ; the covariance parameters and matrix,  $\sigma_u^2$ ,  $\sigma_\varepsilon^2$ ,  $\boldsymbol{\Sigma}^R$ ; and the auxiliary variables,  $\mathbf{a}^R$ . In addition, we parametrise  $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}$  to be a transformation parameter in order to ensure the sampler runs more efficiently. Finally, we write the model in vector notation, which is clearer and much faster in Stan (by making use of more efficient autodifferentiation). The first line specifies the outcome regression model, followed by the prior distributions for the fixed and random effects parameters.

### 2.10.2 Assessment of accuracy

For each simulation setting, we fitted each replicated dataset using both the MFVB and MCMC. The MFVB fits were obtained using the respective Algorithms 3, 4 and 6 with the iterations terminated when the relative increase in  $\log p(\mathbf{y}; q)$  fell below  $10^{-7}$ . The MCMC samples were obtained using Stan (Stan Development Team, 2015) with R (R Development Core Team, 2015) interfacing via the `Rstan` package (Stan Development Team, 2015). In each case, MCMC samples of size 10 000 were generated. The first 5000 values of each sample were discarded as burn-in and the remaining 5000 values were thinned by a factor of 5.

The accuracy of MFVB approximation for a generic parameter  $\theta$  is assessed using the accuracy score defined and justified in Faes *et al.* (2011),

$$\text{accuracy}(q^*(\theta)) = 100 \left( 1 - \frac{1}{2} \int_{-\infty}^{\infty} |q^*(\theta) - p_{\text{MCMC}}(\theta|\mathbf{y})| d\theta \right) \%,$$

This score is derived from the  $L_1$  error, or otherwise known as the *integrated absolute error*. Kernel density estimates, based on the MCMC samples using the binned kernel density estimate `bkde()` function in the R package `KernSmooth` (Wand and Ripley (2006)), were used as a proxy for the exact posterior densities. The default bandwidth is the `dpik()` function in `KernSmooth`, corresponding to the direct plug-in rule. The quality of these proxies dependent upon the MCMC sample sizes and whether or not the posteriors are well-behaved. Nonetheless, it is likely that they are subject to errors inherent in density estimation and bandwidth selection.

Figure 2.7 is an example of summarised MCMC outputs for fitting model (2.11) to the simulated data. Rudimentary diagnostic plots of the MCMC samples given in the first few columns have no discernible pattern, and thus the sampler appears to have mixed well. The kernel density estimated marginal posterior densities of all model parameters, Bayes estimates and 95% credible sets, based on the MCMC samples, are also shown.

Figure 2.8 displays the approximate posterior density functions for the covariance parameters obtained via MFVB and MCMC for a typical realisation from each of the simulation studies. Figure 2.10 displays the side-by-side boxplots of the accuracy scores for the model parameters  $\beta_x$ ,  $\Sigma_{11}^R$ ,  $\Sigma_{12}^R$ ,  $\Sigma_{22}^R$ ,  $\sigma_\varepsilon^2$  and  $f(Q_k)$ ,  $1 \leq k \leq 4$ , where  $Q_k$  are the  $k$ th sample quintiles of the  $s_{ij}$ s. The boxplots show that the majority of the accuracy scores exceed 90% and they rarely drop below 85%. The accuracy drops considerably for the covariance parameters corresponding to the Bernoulli response model. This might be due to the limitations of the Jaakkola and Jordan's approximation as described in Subsection 2.7.1. This deficiency can potentially be remedied through a more elaborate variational approximation - for example, one allows posterior dependence between  $(\boldsymbol{\beta}, \mathbf{u})$  and  $\Sigma^R$ . However, such an elaboration would bring computational costs, which needs to be

traded off against the importance of making inference about the covariance parameters.

In addition, Figure 2.8 indicates that the MFVB and MCMC penalised splines fits are virtually identical, and when overlaid, are indistinguishable from one another, across all response models.

### 2.10.3 Assessment of coverage

Another important type of accuracy assessment is the comparison between the advertised coverage of MFVB credible intervals and the actual coverage. Table 2.2 shows the percentages of true parameter coverage for the approximate 95% credible intervals from the MFVB posterior densities with 0.025 probability mass in each tail, based on 1000 MFVB runs. Across all response models the coverage is generally very good and does not fall below 86% for the majority of the model parameters. For the Bernoulli response model, there is some degradation in coverage accuracy for the covariance parameters.

### 2.10.4 Assessment of speed

To quantify the efficiency gain offered by MFVB, we monitored the time taken for each model to be fitted and the results are summarised in Table 2.1. All computations were performed on a Mac OS X laptop with a 2.6 GHz Intel Core i5 processor and 8 GBytes of random access memory.

As with most speed comparisons, one should be aware of the potential limitations which are: the MFVB and MCMC answers were computed using different programming languages; and both methods have arbitrarily chosen stopping criteria. Despite these caveats, Table 2.1 gives an impression of the relative computational time involved if an “off-the-shelf” MCMC implementation is employed. The results indicate that MFVB is hundred to thousand times faster than MCMC across all models.

Distribution	MCMC	MFVB	Ratio
Gaussian	2768.55 (577.60)	0.52 (0.02)	5324
Student- <i>t</i>	2651.33 (297.09)	13.31 (0.44)	199
Bernoulli	3366.91 (596.45)	29.07 (2.16)	116
Poisson	1061.38 (237.32)	7.60 (1.30)	140

Table 2.1: Average (standard error) elapsed of the computing times in seconds for MCMC and MFVB fitting of the two-level Bayesian semiparametric mixed models to the simulated data.

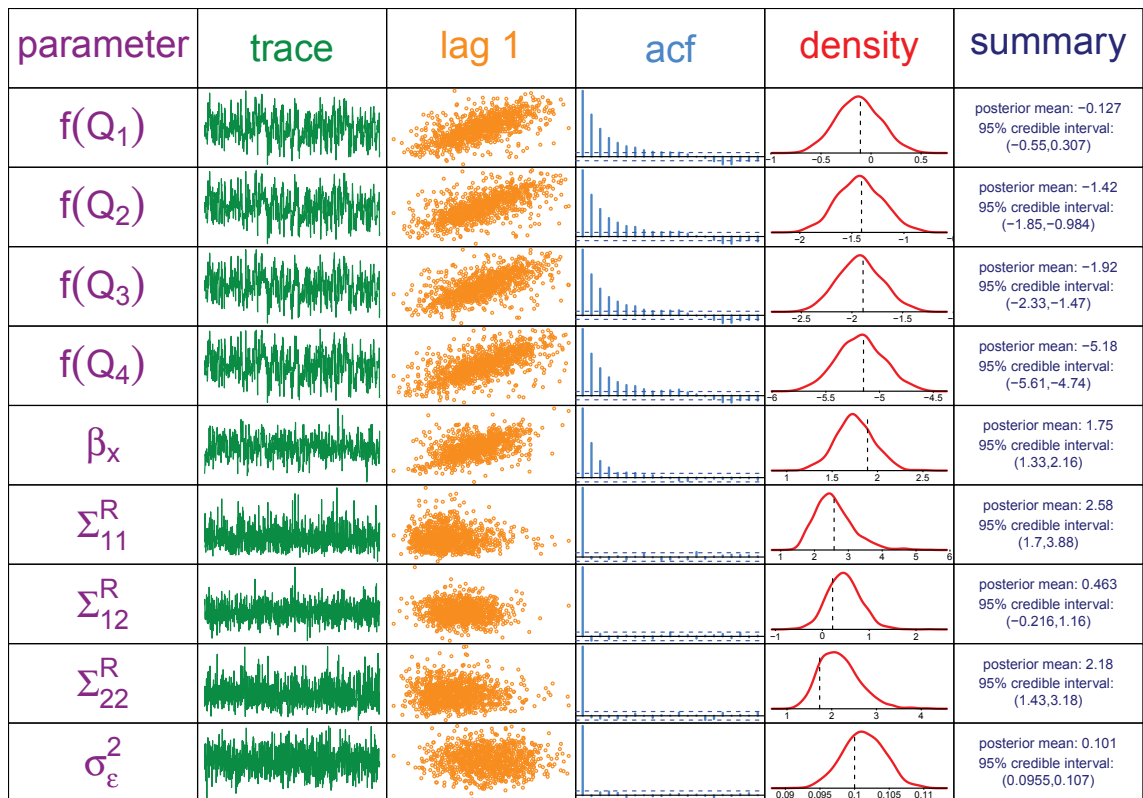


Figure 2.7: Summary of MCMC samples for fitting the two-level Bayesian semiparametric mixed model with Gaussian response to the simulated data. The columns are: parameter name, trace plot of the MCMC sample, plot of sample against its lag-one sample, sample autocorrelation function, kernel density estimate of posterior density function and numerical summaries of posterior density function.

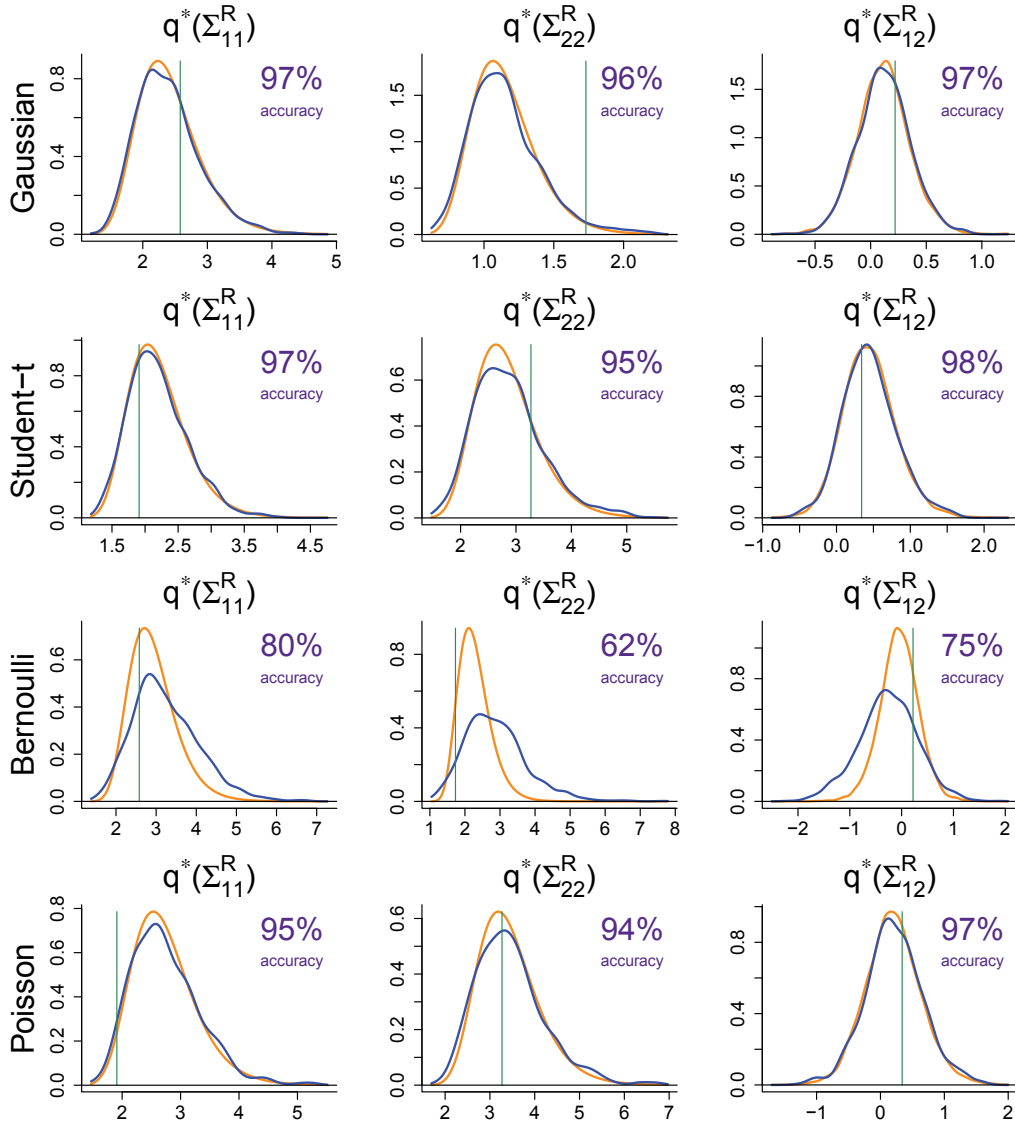


Figure 2.8: Approximate posterior density functions for the covariance parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication from each of the simulation studies described in the text. The green vertical lines represent the true parameter values. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against a MCMC benchmark.



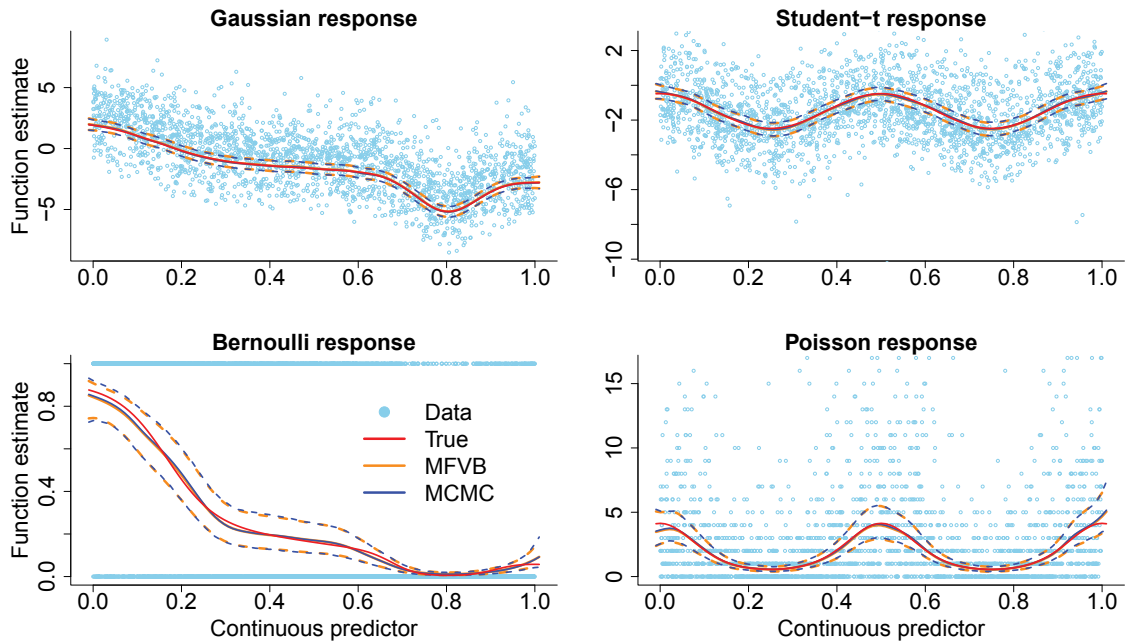


Figure 2.9: Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication of the simulation study described in the text. The red curves represent the true mean functions and the sky blue circles represent the simulated data.

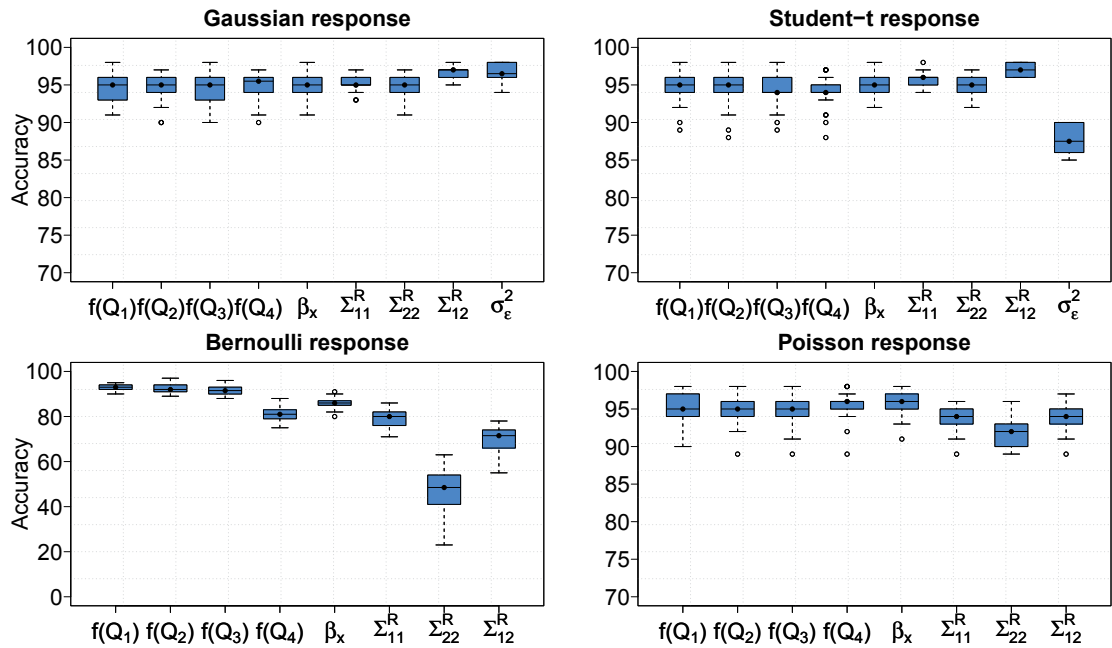


Figure 2.10: Side-by-side boxplots of accuracy scores for MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed models.

## 2.11. CONCLUDING REMARKS

---

Parameter	Gaussian	Student- $t$	Bernoulli	Poisson
$f(Q_1)$	94	98	94	93
$f(Q_2)$	97	97	94	95
$f(Q_3)$	97	93	92	93
$f(Q_4)$	95	96	95	97
$\beta_x$	97	93	92	95
$\Sigma_{11}^R$	95	91	97	94
$\Sigma_{12}^R$	93	91	67	90
$\Sigma_{22}^R$	97	99	68	96
$\sigma_\varepsilon^2$	93	86	—	—

Table 2.2: Percentage coverage of true parameter values by approximate 95% credible sets based on the MFVB approximate posterior density functions. The percentages are based on 1000 replications.

### 2.11 Concluding remarks

Mean field variational Bayes provides a fast and deterministic alternative to MCMC when speed is at a premium. In this chapter we significantly enriched the class of longitudinal and multilevel models which can be handled via MFVB paradigm. The numerical studies in Section 2.10 show that MFVB for various response distributions entails some loss in accuracy for the convenient product restrictions used in our illustrations. Further, particular variational approximations, depending on the problem at hand, can be of hundreds to thousands times faster than the standard MCMC methods. Implementation of these algorithms is often elegant, adding to their practical appeal.

## 2.A Optimal $q$ -densities derivation for Gaussian semiparametric mixed models

We now derive Algorithms 3, 4 and 6 concerning MFVB fitting of the Bayesian semiparametric mixed models. Throughout these appendices, we use the notation “const” to denote additive constants with respect to the function argument. The MFVB calculations rely primarily on the results for the full conditional density functions.

Recall that

$$N \equiv \sum_{i=1}^m n_i, \quad \mathbf{v} \equiv \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \quad \mathbf{C} \equiv [\mathbf{X} \quad \mathbf{Z}], \quad \mathbf{C}_{12} \equiv \text{diag}(b_i)_{1 \leq i \leq N},$$

$$\text{and } \boldsymbol{\Omega} \equiv \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma}^R \end{bmatrix}.$$

**Expressions for  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$**

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \sim N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}),$$

where

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} = \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{y} \quad \text{and}$$

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} = \left( \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes M_{q((\boldsymbol{\Sigma}^R)^{-1})} \end{bmatrix} \right)^{-1}.$$

*Derivation:*

First we note that

$$\begin{aligned} p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) &\propto p(\boldsymbol{\beta}, \mathbf{u} | \mathbf{y}, \sigma_\varepsilon^2, \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) \\ &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \sigma_\varepsilon^2, \boldsymbol{\sigma}_u^2) p(\mathbf{u}^R | \boldsymbol{\Sigma}^R) p(\mathbf{u}^G | \boldsymbol{\sigma}_u^2) p(\boldsymbol{\beta}). \end{aligned}$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) &= -\frac{1}{2\sigma_\varepsilon^2} (\mathbf{y} - \mathbf{C}\mathbf{v})^\top (\mathbf{y} - \mathbf{C}\mathbf{v}) - \frac{1}{2} \mathbf{v}^\top \boldsymbol{\Omega}^{-1} \mathbf{v} + \text{const} \\ &= -\frac{1}{2} \left\{ \mathbf{v}^\top (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C} + \boldsymbol{\Omega}^{-1}) \mathbf{v} - 2\sigma_\varepsilon^{-2} \mathbf{v}^\top \mathbf{C}^\top \mathbf{y} \right\} + \text{const}, \end{aligned}$$

## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

which then leads to the full conditional

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} \sim N \left( (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{y}, (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \right).$$

Taking expectations with respect to all parameters except  $(\boldsymbol{\beta}, \mathbf{u})$  gives

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}) &= E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) \} + \text{const} \\ &= -\frac{1}{2} \left\{ \mathbf{v}^\top \left( \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} + \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \right) \mathbf{v} - 2\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{v}^\top \mathbf{C}^\top \mathbf{y} \right\} + \text{const}, \end{aligned}$$

where

$$\mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \end{bmatrix}.$$

The stated results follow from standard ‘‘completion of square’’ manipulations.

**Expressions for  $q^*(\sigma_\varepsilon^2)$ ,  $B_{q(\sigma_\varepsilon^2)}$  and  $\mu_{q(1/\sigma_\varepsilon^2)}$**

$$q^*(\sigma_\varepsilon^2) \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)} \right),$$

where

$$\begin{aligned} B_{q(\sigma_\varepsilon^2)} &= \frac{1}{2} \left\{ \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top) \right\} + \mu_{q(1/a_\varepsilon)} \\ \text{and } \mu_{q(\sigma_\varepsilon^2)} &= \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)}. \end{aligned}$$

*Derivation:*

First we note that

$$p(\sigma_\varepsilon^2 | \text{rest}) \propto p(\sigma_\varepsilon^2 | \mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, a_\varepsilon) \propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) p(\sigma_\varepsilon^2 | a_\varepsilon).$$

Taking logarithms on both sides gives

$$\log p(\sigma_\varepsilon^2 | \text{rest}) = \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \log(\sigma_\varepsilon^2) - \left( \frac{1}{2} \|\mathbf{y} - \mathbf{C} \mathbf{v}\|^2 + 1/a_\varepsilon \right) / \sigma_\varepsilon^2 + \text{const},$$

## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

which then leads to the full conditional

$$\sigma_\varepsilon^2 | \text{rest} \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), \frac{1}{2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 + 1/a_\varepsilon \right).$$

Taking expectations with respect to all parameters except  $\sigma_\varepsilon^2$  gives

$$\begin{aligned} \log q^*(\sigma_\varepsilon^2) &= E_q \{ \log p(\sigma_\varepsilon^2 | \text{rest}) \} + \text{const} \\ &= \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \log(\sigma_\varepsilon^2) - \left\{ \frac{1}{2} \left( \|\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}\mathbf{C}^\top) \right) \right. \\ &\quad \left. + \mu_{q(1/a_\varepsilon)} \right\} / \sigma_\varepsilon^2 + \text{const}. \end{aligned}$$

The form of  $q^*(\sigma_\varepsilon^2)$  and  $B_{q(\sigma_\varepsilon^2)}$  follows from Result 1.11.6

$$E_q(\|\boldsymbol{\mu}\|^2) = \|E_q(\boldsymbol{\mu})\|^2 + \text{tr}(\text{Cov}_q(\boldsymbol{\mu})).$$

The expression for  $\mu_{q(1/\sigma_\varepsilon^2)}$  follows from elementary manipulations involving inverse-Gamma density functions.

**Expressions for  $q^*(\sigma_{u\ell}^2)$ ,  $B_{q(\sigma_{u\ell}^2)}$  and  $\mu_{q(1/\sigma_{u\ell}^2)}$**

For  $\ell = 1, \dots, L$ ,

$$q^*(\sigma_{u\ell}^2) \sim \text{Inverse-Gamma} \left( \frac{1}{2} (q_\ell^G + 1), B_{q(\sigma_{u\ell}^2)} \right),$$

where

$$B_{q(\sigma_{u\ell}^2)} = \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)}) \right\} + \mu_{q(1/a_{u\ell})} \quad \text{and} \quad \mu_{q(\sigma_{u\ell}^2)} = \frac{1}{2} (q_\ell^G + 1) / B_{q(\sigma_{u\ell}^2)}.$$

*Derivation:*

First we note that

$$p(\sigma_{u\ell}^2 | \text{rest}) \propto p(\sigma_{u\ell}^2 | \mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, a_{u\ell}) \propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_{u\ell}^2) p(\sigma_{u\ell}^2 | a_{u\ell}).$$

Taking logarithms on both sides gives

$$\log p(\sigma_{u\ell}^2 | \text{rest}) = \left\{ -\frac{1}{2} (q_\ell^G + 1) - 1 \right\} \log(\sigma_{u\ell}^2) - \left( \frac{1}{2} \|\mathbf{u}_\ell^G\|^2 + 1/a_{u\ell} \right) / \sigma_{u\ell}^2 + \text{const},$$

## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

which then leads to the full conditional

$$\sigma_{u\ell}^2 | \text{rest} \sim \text{Inverse-Gamma} \left( \frac{1}{2}(q_\ell^G + 1), \frac{1}{2} \|\mathbf{u}_\ell^G\|^2 + 1/a_{u\ell} \right).$$

Taking expectations with respect to all parameters except  $\sigma_{u\ell}^2$  and using Result 1.11.6 gives

$$\begin{aligned} \log q^*(\sigma_{u\ell}^2) &= E_q \{ \log p(\sigma_{u\ell}^2 | \text{rest}) \} + \text{const} \\ &= \left\{ -\frac{1}{2}(q_\ell^G + 1) - 1 \right\} \log(\sigma_{u\ell}^2) - \left\{ \frac{1}{2} \left( \|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)}) \right) \right. \\ &\quad \left. + \mu_{q(1/\sigma_{u\ell}^2)} \right\} / \sigma_{u\ell}^2 + \text{const}. \end{aligned}$$

**Expressions for  $q^*(a_\varepsilon)$ ,  $B_{q(a_\varepsilon)}$  and  $\mu_{q(1/a_\varepsilon)}$**

$$q^*(a_\varepsilon) \sim \text{Inverse-Gamma}(1, B_{q(a_\varepsilon)}),$$

where

$$B_{q(a_\varepsilon)} = \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2} \quad \text{and} \quad \mu_{q(1/a_\varepsilon)} = 1 / \{ \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2} \}.$$

*Derivation:*

First we note that

$$p(a_\varepsilon | \text{rest}) \propto p(a_\varepsilon | \sigma_\varepsilon^2) = p(\sigma_\varepsilon^2 | a_\varepsilon) p(a_\varepsilon).$$

Taking logarithms on both sides gives

$$\log p(a_\varepsilon | \text{rest}) = -2 \log(a_\varepsilon) - (\sigma_\varepsilon^{-2} + A_\varepsilon^{-2}) / a_\varepsilon + \text{const},$$

which then leads to the full conditional

$$a_\varepsilon | \text{rest} \sim \text{Inverse-Gamma} \left( 1, \sigma_\varepsilon^{-2} + A_\varepsilon^{-2} \right).$$

Taking expectations with respect to all parameters except  $a_\varepsilon$  gives

$$\begin{aligned} \log q^*(a_\varepsilon) &= -2 \log(a_\varepsilon) - E_q(\sigma_\varepsilon^{-2} + A_\varepsilon^{-2}) / a_\varepsilon + \text{const} \\ &= (-1 - 1) \log(a_\varepsilon) - (\mu_{q(1/\sigma_\varepsilon^2)} + 1/A_\varepsilon) / a_\varepsilon + \text{const}. \end{aligned}$$

## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

The expressions for  $B_{q(a_\varepsilon)}$  and  $\mu_{q(1/a_\varepsilon)}$  follow immediately.

**Expressions for  $q^*(a_{u\ell}), B_{q(a_{u\ell})}$  and  $\mu_{q(1/a_{u\ell})}$**

$$q^*(a_{u\ell}) \sim \text{Inverse-Gamma}(1, B_{q(a_{u\ell})}),$$

where

$$B_{q(a_{u\ell})} = \mu_{q(1/\sigma_{u\ell}^2)} + A_u^{-2} \quad \text{and} \quad \mu_{q(1/a_{u\ell})} = 1/\{\mu_{q(1/\sigma_{u\ell}^2)} + A_u^{-2}\}.$$

*Derivation:*

The derivation is analogous to that of  $q^*(a_\varepsilon)$  and related quantities.

**Expressions for  $q^*(\Sigma^R), B_{q(\Sigma^R)}$  and  $M_{q((\Sigma^R)^{-1})}$**

$$q^*(\Sigma^R) \sim \text{Inverse-Wishart}(\nu + m + q^R - 1, B_{q(\Sigma^R)}),$$

where

$$\begin{aligned} M_{q((\Sigma^R)^{-1})} &= (\nu + m + q^R - 1) B_{q(\Sigma^R)}^{-1} \quad \text{and} \\ B_{q(\Sigma^R)} &= \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \Sigma_{q(\mathbf{u}_i^R)}) + 2\nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)}). \end{aligned}$$

*Derivation:*

First we note that

$$p(\Sigma^R | \text{rest}) \propto p(\Sigma^R | \mathbf{u}^R, \mathbf{a}^R) = p(\mathbf{u}^R | \Sigma^R) p(\Sigma^R | \mathbf{a}^R).$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(\Sigma^R | \text{rest}) &= -\frac{1}{2} (A_{\Sigma^R} + m + q^R + 1) \log(|\Sigma^R|) \\ &\quad - \frac{1}{2} \text{tr} \left[ \left\{ \sum_{i=1}^m (\mathbf{u}_i^R (\mathbf{u}_i^R)^\top) + 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R) \right\} (\Sigma^R)^{-1} \right] + \text{const}, \end{aligned}$$

## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

which then leads to the full conditional

$$\Sigma^{\mathbf{R}} | \text{rest} \sim \text{Inverse-Wishart} \left( \nu + m + q^{\mathbf{R}} - 1, \sum_{i=1}^m (\mathbf{u}_i^{\mathbf{R}} (\mathbf{u}_i^{\mathbf{R}})^{\top}) + 2\nu \text{diag}(1/a_1^{\mathbf{R}}, \dots, 1/a_{q^{\mathbf{R}}}^{\mathbf{R}}) \right).$$

Taking expectations with respect to all parameters except  $\Sigma^{\mathbf{R}}$  gives

$$\begin{aligned} \log q^*(\Sigma^{\mathbf{R}}) &= -\frac{1}{2} (A_{\Sigma^{\mathbf{R}}} + m + q^{\mathbf{R}} + 1) \log(|\Sigma^{\mathbf{R}}|) \\ &\quad - \frac{1}{2} \text{tr} \left[ E_q \left\{ \sum_{i=1}^m (\mathbf{u}_i^{\mathbf{R}} (\mathbf{u}_i^{\mathbf{R}})^{\top}) + 2\nu \text{diag}(1/a_1^{\mathbf{R}}, \dots, 1/a_{q^{\mathbf{R}}}^{\mathbf{R}}) \right\} (\Sigma^{\mathbf{R}})^{-1} \right] + \text{const} \\ &= -\frac{1}{2} (A_{\Sigma^{\mathbf{R}}} + m + q^{\mathbf{R}} + 1) \log(|\Sigma^{\mathbf{R}}|) \\ &\quad - \frac{1}{2} \text{tr} \left[ \left\{ \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^{\mathbf{R}})} \boldsymbol{\mu}_{q(\mathbf{u}_i^{\mathbf{R}})}^{\top} + \Sigma_{q(\mathbf{u}_i^{\mathbf{R}})}) \right. \right. \\ &\quad \quad \left. \left. + 2\nu \text{diag}(\mu_{q(1/a_1^{\mathbf{R}})}, \dots, \mu_{q(1/a_{q^{\mathbf{R}}}^{\mathbf{R}})}) \right\} (\Sigma^{\mathbf{R}})^{-1} \right] + \text{const}. \end{aligned}$$

The expressions for  $B_{q(\Sigma^{\mathbf{R}})}$  and  $M_{q((\Sigma^{\mathbf{R}})^{-1})}$  follow immediately.

**Expressions for  $q^*(a_r^{\mathbf{R}})$ ,  $B_{q(a_r^{\mathbf{R}})}$  and  $\mu_{q(a_r^{\mathbf{R}})}$**

For  $r = 1, \dots, q^{\mathbf{R}}$ ,

$$q^*(a_r^{\mathbf{R}}) \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\frac{1}{2}(\nu + q^{\mathbf{R}}), B_{q(a_r^{\mathbf{R}})}),$$

where

$$B_{q(a_r^{\mathbf{R}})} = \nu (M_{q((\Sigma^{\mathbf{R}})^{-1})})_{rr} + A_{\mathbf{R}r}^{-2} \quad \text{and} \quad \mu_{q(a_r^{\mathbf{R}})} = \frac{1}{2}(\nu + q^{\mathbf{R}}) / B_{q(a_r^{\mathbf{R}})}.$$

*Derivation:*

First we note that

$$p(a_r^{\mathbf{R}} | \text{rest}) \propto p(a_r^{\mathbf{R}} | \Sigma^{\mathbf{R}}) = p(\Sigma^{\mathbf{R}} | a_r^{\mathbf{R}}) p(a_r^{\mathbf{R}}).$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(a_r^{\mathbf{R}} | \text{rest}) &= -\left\{ \frac{1}{2}(\nu + q^{\mathbf{R}} - 1) \right\} (\log(2\nu) + \log(a_r^{\mathbf{R}})) - \left\{ \nu (\Sigma^{\mathbf{R}})^{-1} \right\}_{rr} / a_r^{\mathbf{R}} \\ &\quad - (a_r^{\mathbf{R}})^{-1} A_{\mathbf{R}r}^{-2}, \end{aligned}$$



## 2.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

which then leads to the full conditional

$$a_r^R | \text{rest} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}(\nu + q^R), \nu(\boldsymbol{\Sigma}^R)_{rr}^{-1} + A_{Rr}^{-2}\right).$$

Taking expectations with respect to all parameters except  $a_r^R$  gives

$$\begin{aligned} \log q^*(a_r^R) &= -\left\{\frac{1}{2}(\nu + q^R - 1)\right\} (\log(2\nu) + \log(a_r^R)) - \nu(\mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})})_{rr} / a_r^R \\ &\quad - (a_r^R)^{-1} A_{Rr}^{-2}. \end{aligned}$$

The expressions for  $B_{q(a_r^R)}$  and  $\mu_{q(a_r^R)}$  follow immediately.

### 2.A.1 Derivation of the marginal log-likelihood lower bound

The expression for the marginal log-likelihood lower bound is

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^R (\nu + q^R - 1) \log(2\nu) - \frac{1}{2} \log(2\pi) \sum_{i=1}^m n_i - \left(\frac{1}{2} q^R + L + 1\right) \log(\pi) \\ &\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) \\ &\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| - \log(\mathcal{C}_{q^R, \nu + q^R - 1}) + \log(\mathcal{C}_{q^R, \nu + m + q^R - 1}) \\ &\quad - \frac{1}{2} (\nu + m + q^R - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}| \\ &\quad + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^G + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\ &\quad + \log \Gamma \left\{ \frac{1}{2} (\sum_{i=1}^m n_i + 1) \right\} - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_\varepsilon^2)}) \\ &\quad - \sum_{r=1}^{q^R} \log(A_{Rr}) + q^R \log \Gamma \left\{ \frac{1}{2} (\nu + q^R) \right\} \\ &\quad - \sum_{\ell=1}^L \log(A_{u\ell}) - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\ &\quad + \sum_{r=1}^{q^R} \nu(\mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})})_{rr} \mu_{q(1/a_r^R)} - \frac{1}{2} (\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) \\ &\quad - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)}. \end{aligned}$$

*Derivation:*

The logarithm of the lower bound on the marginal likelihood is given by

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon) \right. \\ &\quad \left. - \log q^*(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \boldsymbol{\Sigma}^R, \sigma_u^2, \sigma_\varepsilon^2, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon) \right\}. \end{aligned}$$

According to the product restriction (2.12), our expression for the lower bound of the marginal log-likelihood becomes:

$$\log \underline{p}(\mathbf{y}; q) = E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2) \right\}$$

2.A. OPTIMAL  $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

$$\begin{aligned}
& + E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) \} \\
& + E_q \{ \log p(\boldsymbol{\Sigma}^R | \mathbf{a}^R) - \log q^*(\boldsymbol{\Sigma}^R) \} + E_q \{ p(\boldsymbol{\sigma}_u^2 | \mathbf{a}_u) - \log q^*(\boldsymbol{\sigma}_u^2) \} \\
& + E_q \{ \log p(\mathbf{a}^R) - \log q^*(\mathbf{a}^R) \} + E_q \{ \log p(\mathbf{a}_u) - \log q^*(\mathbf{a}_u) \} \\
& + E_q \{ p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \} + E_q \{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \}.
\end{aligned}$$

First note that

$$\log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2) = \frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2.$$

Taking expectations gives

$$\begin{aligned}
E_q \{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2) \} & = \frac{N}{2} \log(2\pi) - \frac{N}{2} E_q \{ \log(\sigma_\varepsilon^2) \} \\
& - \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \|\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\}.
\end{aligned}$$

Using the well-established result concerning entropy of a multivariate random normal vector results in

$$\begin{aligned}
E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) \} & = \\
& - \frac{p}{2} \log(2\pi) - \frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{2\sigma_\beta^2} (\|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})})) \\
& - \frac{m}{2} \log(2\pi) - \frac{m}{2} E_q \{ \log(|\boldsymbol{\Sigma}^R|) \} - \frac{1}{2} (\mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})}) (\|\boldsymbol{\mu}_{q(\mathbf{u}^R)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}^R)})) \\
& - \frac{1}{2} \sum_{\ell=1}^L \left\{ q_\ell^G \log(2\pi) - q_\ell^G E_q \{ \log(|\sigma_{u\ell}^2|) \} - \frac{1}{2} \mu_{q(1/\sigma_{u\ell}^2)} (\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)})) \right\} \\
& - \frac{1}{2} \log(|\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|) - \frac{1}{2} (m + p + \sum_{\ell=1}^L q_\ell^G) \log(2\pi) - \frac{1}{2} (m + p + \sum_{\ell=1}^L q_\ell^G).
\end{aligned}$$

The next term is

$$\begin{aligned}
\log p(\boldsymbol{\Sigma}^R | \mathbf{a}^R) - \log q^*(\boldsymbol{\Sigma}^R) & = \\
& \log(\mathcal{C}_{q^R, v+m-1}) + \log(\mathcal{C}_{q^R, v+m+q^R-1}) + \frac{1}{2} (v + m - 1) \log(|\mathbf{B}_{\boldsymbol{\Sigma}^R}|) \\
& - \frac{1}{2} (v + 2m - 1) \log(|\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}|) + \frac{1}{2} (v + m + q^R) \log(|\boldsymbol{\Sigma}^R|) \\
& + \frac{1}{2} (v + 2m + q^R) \log(|\boldsymbol{\Sigma}^R|) - \frac{1}{2} \text{tr}(\mathbf{B}_{\boldsymbol{\Sigma}^R} (\boldsymbol{\Sigma}^R)^{-1}) + \text{tr}(\mathbf{B}_{q(\boldsymbol{\Sigma}^R)} (\boldsymbol{\Sigma}^R)^{-1}).
\end{aligned}$$

Taking expectations gives

$$\begin{aligned}
E_q \{ \log p(\boldsymbol{\Sigma}^R | \mathbf{a}^R) - \log q^*(\boldsymbol{\Sigma}^R) \} & = \\
& \log(\mathcal{C}_{q^R, v+m-1}) + \log(\mathcal{C}_{q^R, v+2m-1}) + \frac{1}{2} (v + m - 1) \log(|\mathbf{B}_{\boldsymbol{\Sigma}^R}|) \\
& - \frac{1}{2} (v + 2m - 1) \log(|\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}|) + \frac{m}{2} E_q \{ \log(|\boldsymbol{\Sigma}^R|) \}
\end{aligned}$$

2.A. OPTIMAL  $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

$$+ \frac{1}{2} \text{tr} \left\{ \left( \mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})} - \mathbf{B}_{\boldsymbol{\Sigma}^{\text{R}}} \right) \mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \right\},$$

where

$$\begin{aligned} \frac{1}{2} \text{tr} \left\{ \left( \mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})} - \mathbf{B}_{\boldsymbol{\Sigma}^{\text{R}}} \right) \mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \right\} &= \frac{1}{2} \text{tr} \left\{ \sum_{i=1}^m \left( \boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})} \boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})}^{\text{T}} + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^{\text{R}})} \right) \mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \right\} \\ &= \frac{1}{2} \left( \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \right) \left( \boldsymbol{\mu}_{q(\mathbf{u}^{\text{R}})}^{\text{T}} \boldsymbol{\mu}_{q(\mathbf{u}^{\text{R}})} + \boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{R}})} \right). \end{aligned}$$

The next term is

$$\begin{aligned} \log p(\sigma_{u\ell}^2 | a_{u\ell}) - \log q^*(\sigma_{u\ell}^2) &= \\ \log \left\{ \frac{a_{u\ell}^{-0.5}}{\Gamma(0.5)} (\sigma_{u\ell}^2)^{-1.5} \exp \left( -\frac{1}{a_{u\ell} \sigma_{u\ell}^2} \right) \right\} & \\ - \log \left\{ \frac{(B_{q(\sigma_{u\ell}^2)})^{0.5(q_\ell^{\text{G}}+1)}}{\Gamma \left\{ \frac{1}{2} (q_\ell^{\text{G}} + 1) \right\}} - (\sigma_{u\ell}^2)^{\left\{ -\frac{1}{2} (q_\ell^{\text{G}}+1) - 1 \right\}} \exp \left( -\frac{B_{q(\sigma_{u\ell}^2)}}{\sigma_{u\ell}^2} \right) \right\}. & \end{aligned}$$

Taking expectations gives

$$\begin{aligned} E_q \{ \log p(\sigma_{u\ell}^2 | a_{u\ell}) - \log q^*(\sigma_{u\ell}^2) \} &= \\ -\frac{1}{2} E_q \{ \log(a_{u\ell}) \} - \frac{1}{2} \log(\pi) + \log \Gamma \left\{ \frac{1}{2} (q_\ell^{\text{G}} + 1) \right\} & \\ + \frac{1}{2} q_\ell^{\text{G}} E_q \{ \log(\sigma_{u\ell}^2) \} + \left( B_{q(\sigma_{u\ell}^2)} - \mu_{q(1/a_{u\ell})} \right) \mu_{q(1/\sigma_{u\ell}^2)} & \\ - \frac{1}{2} (q_\ell^{\text{G}} + 1) \log(B_{q(\sigma_{u\ell}^2)}), & \end{aligned}$$

where

$$(B_{q(\sigma_{u\ell}^2)} - \mu_{q(1/a_{u\ell})}) \mu_{q(1/\sigma_{u\ell}^2)} = \frac{1}{2} \mu_{q(1/\sigma_{u\ell}^2)} (\|\mathbf{u}_\ell^{\text{G}}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^{\text{G}})})).$$

Similarly, for

$$\begin{aligned} E_q \{ \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \} &= -\frac{1}{2} E_q \{ \log(a_\varepsilon) \} - \frac{1}{2} \log(\pi) + \log \Gamma \left\{ \frac{1}{2} (N + 1) \right\} \\ &+ \frac{N}{2} E_q \{ \log(\sigma_\varepsilon^2) \} + (B_{q(\sigma_\varepsilon^2)} - \mu_{q(1/a_\varepsilon)}) \mu_{q(1/\sigma_\varepsilon^2)} \\ &- \frac{1}{2} (N + 1) \log \{ B_{q(\sigma_\varepsilon^2)} \}. \end{aligned}$$

The next term is

$$\log p(a_r^{\text{R}}) - \log q^*(a_r^{\text{R}}) = \log \left\{ \prod_{r=1}^{q^{\text{R}}} \frac{A_{\text{R}r}^{-1}}{\Gamma(0.5)} (a_r^{\text{R}})^{-1.5} \exp \left( \frac{-1}{a_r^{\text{R}} A_{\text{R}r}} \right) \right\}$$

2.A. OPTIMAL  $q$ -DENSITIES DERIVATION FOR GAUSSIAN SEMIPARAMETRIC MIXED MODELS

---

$$-\log \left\{ \prod_{r=1}^{q^R} \frac{(B_{q(a_r^R)})^{0.5(\nu+2)}}{\Gamma\left\{\frac{1}{2}(\nu+2)\right\}} (a_r^R)^{\left\{-\frac{1}{2}(\nu+1)-1\right\}} \exp\left(-\frac{B_{q(a_r^R)}}{a_r^R}\right) \right\}.$$

Taking expectations gives

$$\begin{aligned} E_q\{\log p(a_r^R) - \log q^*(a_r^R)\} = & \\ & \frac{1}{2} \sum_{r=1}^{q^R} \log A_{Rr} - \log(\pi) + 2 \log \Gamma\left\{\frac{1}{2}(\nu+2)\right\} + \frac{1}{2}(\nu+q^R-1) \sum_{r=1}^{q^R} E_q\{\log(a_r^R)\} \\ & + \sum_{r=1}^{q^R} \left\{ (B_{q(a_r^R)} - 1)/A_{Rr} \right\} \mu_{q(1/a_r^R)} - \frac{1}{2}(\nu+q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}), \end{aligned}$$

where

$$\sum_{r=1}^{q^R} \left\{ (B_{q(a_r^R)} - 1)/A_{Rr} \right\} \mu_{q(1/a_r^R)} = \sum_{r=1}^{q^R} \nu (M_{q((\Sigma^R)^{-1})})_{rr} \mu_{q(1/a_r^R)}.$$

The next term is

$$\begin{aligned} \log p(a_{u\ell}) - \log q^*(a_{u\ell}) = & \log \left\{ \frac{(A_{u\ell})^{-1}}{\Gamma(0.5)} (a_{u\ell})^{-1.5} \exp\left(-\frac{1}{a_{u\ell} A_{u\ell}^2}\right) \right\} \\ & - \log \left\{ \frac{B_{q(a_{u\ell})}}{\Gamma(1)} (a_{u\ell})^{-2} \exp\left(-\frac{B_{q(a_{u\ell})}}{a_{u\ell}}\right) \right\}. \end{aligned}$$

Taking expectations gives

$$\begin{aligned} E_q\{\log p(a_{u\ell}) - \log q^*(a_{u\ell})\} = & -\frac{1}{2} \log(A_{u\ell}^2) - \log(B_{q(a_{u\ell})}) - \frac{1}{2} \log(\pi) + \frac{1}{2} E_q\{\log(a_{u\ell})\} \\ & + (B_{q(a_{u\ell})} - 1/A_{u\ell}^2) \mu_{q(1/a_{u\ell})}, \end{aligned}$$

where

$$(B_{q(a_{u\ell})} - 1/A_{u\ell}^2) \mu_{q(1/a_{u\ell})} = \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)}.$$

The stated results then follows from cancellations lead by these quantities, as well as more obvious ones.

## 2.B Optimal $q$ -densities derivation for Student- $t$ semiparametric mixed models

Expressions for  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \sim \text{Normal}(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}),$$

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y} \quad \text{and} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \left( \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u^\ell}^2)} \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \end{bmatrix}^{-1} \right). \end{aligned}$$

*Derivation:*

First we note that

$$\begin{aligned} p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) &\propto p(\boldsymbol{\beta}, \mathbf{u} | \mathbf{y}, \sigma_\varepsilon^2, \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) \\ &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2) p(\mathbf{u}^G | \boldsymbol{\sigma}_u^2) p(\mathbf{u}^R | \boldsymbol{\Sigma}^R) p(\boldsymbol{\beta}). \end{aligned}$$

Taking logarithms on both sides gives

$$\log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) = -\frac{1}{2} \left\{ \mathbf{v}^\top (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \boldsymbol{\Omega}^{-1}) \mathbf{v} - 2\sigma_\varepsilon^{-2} \mathbf{v}^\top \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y} \right\} + \text{const},$$

which leads to the full conditional

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} \sim N \left( (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y}, (\sigma_\varepsilon^{-2} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \right).$$

Taking expectations with respect to all parameters except  $(\boldsymbol{\beta}, \mathbf{u})$  gives

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}) &= E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) \} + \text{const} \\ &= -\frac{1}{2} \left\{ \mathbf{v}^\top \left( \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \right) \mathbf{v} - 2\sigma_\varepsilon^{-2} \mathbf{v}^\top \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y} \right\} + \text{const}. \end{aligned}$$

The stated results follows from standard ‘‘completion of square’’ manipulations.

The remaining  $q^*$  densities involve straightforward adaptation of the derivations given in the Gaussian response case.

### 2.B.1 Derivation of the marginal log-likelihood lower bound

The expression for the marginal log-likelihood lower bound is

$$\begin{aligned}
\underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^{\text{R}}(\nu + q^{\text{R}} - 1) \log(2\nu) - \frac{1}{2} \log(2\pi) \sum_{i=1}^m n_i - \left(\frac{1}{2} q^{\text{R}} + L + 1\right) \log(\pi) \\
&\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^{\text{G}} + p + m \right) \\
&\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| - \log(\mathcal{C}_{q^{\text{R}}, \nu + q^{\text{R}} - 1}) + \log(\mathcal{C}_{q^{\text{R}}, \nu + m + q^{\text{R}} - 1}) \\
&\quad - \frac{1}{2} (\nu + m + q^{\text{R}} - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})}| \\
&\quad + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^{\text{G}} + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^{\text{G}} + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\
&\quad + \log \Gamma \left\{ \frac{1}{2} (\sum_{i=1}^m n_i + 1) \right\} - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_\varepsilon^2)}) \\
&\quad - \sum_{r=1}^{q^{\text{R}}} \log(A_{\text{R}r}) + q^{\text{R}} \log \Gamma \left\{ \frac{1}{2} (\nu + q^{\text{R}}) \right\} \\
&\quad - \sum_{\ell=1}^L \log(A_{u\ell}) - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\
&\quad + \sum_{r=1}^{q^{\text{R}}} \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})_{rr}}) \mu_{q(1/a_r^{\text{R}})} - \frac{1}{2} (\nu + q^{\text{R}}) \sum_{r=1}^{q^{\text{R}}} \log(B_{q(a_r^{\text{R}})}) \\
&\quad - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)} \\
&\quad - \frac{1}{2} \mu_{q(\nu+1)} \sum_{i=1}^N \log B_{q(b_i)} + \log \Gamma \left\{ \frac{1}{2} (\mu_{q(\nu)} + 1) \right\} + \sum_{i=1}^N B_{q(b_i)} \mu_{q(1/b_i)} \\
&\quad - \log(\nu_{\max} - \nu_{\min}) - \frac{1}{2} \mu_{q(\nu)} \sum_{i=1}^N \mu_{q(\log(b_i))} + \log \mathcal{F}(0, \sum_{i=1}^m n_i, C_1, \nu_{\max}, \nu_{\min}),
\end{aligned}$$

*Derivation:*

The logarithm of the lower bound on the marginal likelihood is given by

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mathbf{b}, v) \right. \\
&\quad \left. - \log q^*(\boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mathbf{b}, v) \right\}.
\end{aligned}$$

According to the product restriction (2.17), our expression for the lower bound of the marginal log-likelihood becomes:

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \sigma_\varepsilon^2, \mathbf{b}) \right\} \\
&\quad + E_q \left\{ \log p(\boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}} | \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}) \right\} \\
&\quad + E_q \left\{ \log p(\boldsymbol{\Sigma}^{\text{R}} | \mathbf{a}^{\text{R}}) - \log q^*(\boldsymbol{\Sigma}^{\text{R}}) \right\} + E_q \left\{ p(\boldsymbol{\sigma}_u^2 | \mathbf{a}_u) - \log q^*(\boldsymbol{\sigma}_u^2) \right\} \\
&\quad + E_q \left\{ \log p(\mathbf{a}^{\text{R}}) - \log q^*(\mathbf{a}^{\text{R}}) \right\} + E_q \left\{ \log p(\mathbf{a}_u) - \log q^*(\mathbf{a}_u) \right\} \\
&\quad + E_q \left\{ p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \right\} + E_q \left\{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \right\} \\
&\quad + E_q \left\{ \log p(\mathbf{b} | v) - \log q^*(\mathbf{b}) \right\} + E_q \left\{ \log p(v) - q^*(v) \right\}.
\end{aligned}$$

We are only interested in the first and last terms (highlighted in darker colour) since they are the only differences between the Gaussian and Student- $t$  response cases.

2.B. OPTIMAL  $q$ -DENSITIES DERIVATION FOR STUDENT- $t$  SEMIPARAMETRIC MIXED MODELS

---

First note that

$$\begin{aligned} \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2, \mathbf{b}) &= \frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2} \sum_{i=1}^N \log(1/b_i) \\ &\quad - \frac{1}{2} \left\{ (\mathbf{y} - \mathbf{C}\mathbf{v})^\top (\sigma_\varepsilon^2 \mathbf{C}_{12})^{-1} (\mathbf{y} - \mathbf{C}\mathbf{v}) \right\}. \end{aligned}$$

Taking expectations gives

$$\begin{aligned} E_q \{ \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2, \mathbf{b}) \} &= \frac{N}{2} \log(2\pi) - \frac{N}{2} E_q \{ \log(\sigma_\varepsilon^2) \} - \frac{1}{2} \sum_{i=1}^N \log(1/\mu_q(b_i)) \\ &\quad - \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right. \\ &\quad \left. + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C}) \right\}. \end{aligned}$$

The next term is

$$\begin{aligned} \log p(\mathbf{b}|v) - \log q^*(\mathbf{b}) &= N \left\{ \frac{v}{2} \log(v/2) - \log \Gamma(v/2) \right\} - \frac{v}{2} \sum_{i=1}^N \{ \log(b_i) + (1/b_i) \} \\ &\quad - \frac{1}{2} (\mu_{q(v)} + 1) \sum_{i=1}^N \log B_{q(b_i)} + N \log \Gamma \left\{ \frac{1}{2} \mu_{q(v)} + 1 \right\} \\ &\quad + \frac{1}{2} (\mu_{q(v)} + 3) \sum_{i=1}^N \log(b_i) - \sum_{i=1}^N B_{q(b_i)} / b_i. \end{aligned}$$

Taking expectations gives

$$\begin{aligned} E_q \{ \log p(\mathbf{b}|v) - \log q^*(\mathbf{b}) \} &= E_q \left\{ N \left( \frac{v}{2} \log(v/2) - \log \Gamma(v/2) \right) \right\} \\ &\quad - \frac{1}{2} (\mu_{q(v)} + 1) \sum_{i=1}^N \log B_{q(b_i)} + N \log \Gamma \left\{ \frac{1}{2} (\mu_{q(v)} + 1) \right\} \\ &\quad - \frac{1}{2} \sum_{i=1}^N \mu_{q(\log(b_i))} + \sum_{i=1}^N (B_{q(b_i)} - \frac{1}{2} \mu_{q(v)}) \mu_{q(1/b_i)}. \end{aligned}$$

The next term is

$$\begin{aligned} \log p(v) - \log q^*(v) &= -\log(v_{\max} - v_{\min}) - N \left( \frac{v}{2} \log(v/2) - \log \Gamma(v/2) \right) \\ &\quad + \frac{v}{2} \sum_{i=1}^N \{ \log(b_i) + (1/b_i) \} + \log \mathcal{F}(0, N, \mathbf{C}_1, v_{\max}, v_{\min}). \end{aligned}$$

## 2.B. OPTIMAL $q$ -DENSITIES DERIVATION FOR STUDENT- $t$ SEMIPARAMETRIC MIXED MODELS

---

Taking expectations gives

$$\begin{aligned} E_q \{ \log p(v) - \log q^*(v) \} &= -\log(v_{\max} - v_{\min}) - E_q \left( N \left( \frac{v}{2} \log(v/2) - \log \Gamma(v/2) \right) \right) \\ &\quad + \frac{1}{2} \mu_{q(v)} \sum_{i=1}^N \{ \mu_{q(\log(b_i))} + \mu_{q(1/b_i)} \} + \log \mathcal{F}(0, N, \mathbf{C}_1, v_{\max}, v_{\min}). \end{aligned}$$

**Expressions for  $q^*(\sigma_\varepsilon^2)$ ,  $B_{q(\sigma_\varepsilon^2)}$  and  $\mu_{q(1/\sigma_\varepsilon^2)}$**

$$q^*(\sigma_\varepsilon^2) \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)} \right),$$

where

$$\begin{aligned} \mu_{q(\sigma_\varepsilon^2)} &= \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)} \quad \text{and} \\ B_{q(\sigma_\varepsilon^2)} &= \frac{1}{2} \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C}) \right\} + \mu_{q(1/a_\varepsilon)}. \end{aligned}$$

*Derivation:*

First we note that

$$p(\sigma_\varepsilon^2 | \text{rest}) \propto p(\sigma_\varepsilon^2 | \mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, a_\varepsilon) \propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) p(\sigma_\varepsilon^2 | a_\varepsilon).$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(\sigma_\varepsilon^2 | \text{rest}) &= \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \\ &\quad - \log(\sigma_\varepsilon^2) \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{C} \mathbf{v})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C} \mathbf{v}) + 1/a_\varepsilon \right\} / \sigma_\varepsilon^2 + \text{const}, \end{aligned}$$

which then leads to the full conditional

$$\sigma_\varepsilon^2 | \text{rest} \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), \frac{1}{2} (\mathbf{y} - \mathbf{C} \mathbf{v})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C} \mathbf{v}) + a_\varepsilon^{-1} \right).$$

Taking expectations with respect to all parameters except  $\sigma_\varepsilon^2$  gives

$$\begin{aligned} \log q^*(\sigma_\varepsilon^2) &= E_q \{ \log p(\sigma_\varepsilon^2 | \text{rest}) \} + \text{const} \\ &= \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \log(\sigma_\varepsilon^2) - \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right. \\ &\quad \left. + \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C}) + \mu_{q(1/a_\varepsilon)} \right\} / \sigma_\varepsilon^2 + \text{const}. \end{aligned}$$

The expression for  $\mu_{q(1/\sigma_\varepsilon^2)}$  follows from elementary manipulations involving inverse-Gamma density functions.



2.B. OPTIMAL  $q$ -DENSITIES DERIVATION FOR STUDENT- $t$  SEMIPARAMETRIC MIXED MODELS

---

**Expressions for  $q^*(b_i)$ ,  $B_{q(b_i)}$  and  $\mu_{q(b_i)}$**

$$q^*(\mathbf{b}) = \prod_{i=1}^N q^*(b_i); \quad q^*(b_i) \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}(\mu_{q(\nu)} + 1), B_{q(b_i)}\right)$$

where

$$B_{q(b_i)} = \frac{1}{2} \left\{ \mu_{q(\nu)} + (\mathbf{y} - \mathbf{C}\mathbf{v})_i^2 / \sigma_\varepsilon^2 \right\} \quad \text{and} \quad \mu_{q(b_i)} = \frac{1}{2}(\mu_{q(\nu)} + 1) / B_{q(b_i)}.$$

*Derivation:*

First we note that

$$p(\mathbf{b}|\text{rest}) \propto p(\mathbf{b}|\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2, \nu) \propto p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \mathbf{b}, a_\varepsilon) p(\mathbf{b}|\nu).$$

Taking logarithms on both sides gives

$$\log p(\mathbf{b}|\text{rest}) = \left\{ -\frac{1}{2}(\nu + 1) - 1 \right\} \sum_{i=1}^N \left[ \log(b_i) - \frac{1}{2b_i} \left\{ \nu + (\mathbf{y} - \mathbf{C}\mathbf{v})_i^2 / \sigma_\varepsilon^2 \right\} \right] + \text{const},$$

which then leads to the full conditional

$$b_i|\text{rest} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}(\nu + 1), \frac{1}{2} \left\{ \nu + (\mathbf{y} - \mathbf{C}\mathbf{v})_i^2 / \sigma_\varepsilon^2 \right\} \right).$$

Taking expectations with respect to all parameters except  $\mathbf{b}$  gives

$$\begin{aligned} \log q^*(\mathbf{b}) &= \left\{ -\frac{1}{2}(\mu_{q(\nu)} + 1) - 1 \right\} \sum_{i=1}^N \left[ \log(b_i) - \frac{1}{2b_i} \left\{ \mu_{q(\nu)} + \left( (\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})_i^2 \right. \right. \right. \\ &\quad \left. \left. \left. + (\mathbf{C}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}\mathbf{C}^\top)_{ii} / \mu_{q(\sigma_\varepsilon^2)} \right) \right\} \right] + \text{const}. \end{aligned}$$

The expressions for  $B_{q(b_i)}$  and  $\mu_{q(b_i)}$  follow immediately.

**Expressions for  $p(\nu|\text{rest})$  and  $q^*(\nu)$**

The optimal density  $q^*(\nu)$  is not of a standard form:

$$\begin{aligned} p(\nu|\text{rest}) &\propto p(\nu|\mathbf{b}) = p(\mathbf{b}|\nu) p(\nu) \\ &= \frac{1}{(\nu_{\max} - \nu_{\min})} \prod_{i=1}^N \left\{ \frac{(0.5\nu)^{0.5\nu}}{\Gamma(0.5\nu)} (b_i)^{-\left(\frac{1}{2}\nu+1\right)} \exp\left(-\frac{\nu}{2b_i}\right) \right\} \end{aligned}$$

2.C. OPTIMAL  $q$ -DENSITIES DERIVATION FOR BERNOULLI SEMIPARAMETRIC MIXED MODELS

---

$$\begin{aligned}\log p(\nu|\text{rest}) &= N \left\{ \frac{\nu}{2} \log(\nu/2) - \log \Gamma(\nu/2) \right\} - \frac{\nu}{2} \sum_{i=1}^N (\log(b_i) + b_i^{-1}) + \text{const} \\ \log q^*(\nu) &= N \left\{ \frac{\nu}{2} \log(\nu/2) - \log \Gamma(\nu/2) \right\} - \frac{\nu}{2} \sum_{i=1}^N (\log(\mu_{q(b_i)}) + \mu_{q(1/b_i)}) + \text{const}.\end{aligned}$$

The remaining  $q^*$  densities involve straightforward adaptation of the derivations given in the Gaussian response case.

## 2.C Optimal $q$ -densities derivation for Bernoulli semiparametric mixed models

Expressions for  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}$

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \sim \text{Normal}(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}),$$

where

$$\begin{aligned}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \mathbf{C}^\top (\mathbf{y} - \tfrac{1}{2} \mathbf{1}) \quad \text{and} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} &= \left( 2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \end{bmatrix} \right)^{-1}.\end{aligned}$$

*Derivation:*

First note that

$$\begin{aligned}p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) &= \exp \left\{ \mathbf{y}^\top \mathbf{C} \mathbf{v} - \left( \mathbf{v}^\top \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} \mathbf{v} - \tfrac{1}{2} \mathbf{1}^\top \mathbf{C} \mathbf{v} + \tfrac{1}{2} \mathbf{1}^\top \boldsymbol{\psi}(\boldsymbol{\xi}) \right) \right. \\ &\quad \left. - \tfrac{1}{2} \mathbf{v}^\top \boldsymbol{\Omega}^{-1} \mathbf{v} \right\} + \text{const}.\end{aligned}$$

Taking logarithms on both sides gives

$$\log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) = -\tfrac{1}{2} \left\{ \mathbf{v}^\top \left( 2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \boldsymbol{\Omega}^{-1} \right) \mathbf{v} - 2\mathbf{C}^\top (\mathbf{y} - \tfrac{1}{2} \mathbf{1}) \mathbf{v}^\top \right\} + \text{const},$$

which then leads to the full conditional

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Bigg|_{\text{rest}} \sim N \left( (2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1} \mathbf{C}^\top (\mathbf{y} - \tfrac{1}{2} \mathbf{1}), \right.$$

$$(2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\}\mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1}).$$

Taking expectations with respect to all parameters except  $(\boldsymbol{\beta}, \mathbf{u})$  gives

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) &= E_q \{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) | \text{rest} \} + \text{const} \\ &= -\frac{1}{2} \left\{ \mathbf{v}^\top \left( 2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\}\mathbf{C} + \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \right) \mathbf{v} - 2\mathbf{C}^\top \left( \mathbf{y} - \frac{1}{2}\mathbf{1} \right) \mathbf{v}^\top \right\} + \text{const}, \end{aligned}$$

where the optimal update for the additional variational vector  $\boldsymbol{\xi}$  is

$$\boldsymbol{\xi} \leftarrow \sqrt{\text{diagonal} \left\{ \mathbf{C} \left( \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^\top \right) \mathbf{C}^\top \right\}}.$$

This follows the Expectation-Maximisation arguments presented in Jaakkola and Jordan (2000). The remaining  $q^*$  densities involve straightforward adaptation of the derivations given in the Gaussian response case.

### 2.C.1 Derivation of the marginal log-likelihood lower bound

The expression for the marginal log-likelihood lower bound is

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} q^{\text{R}} (\nu + q^{\text{R}} - 1) \log(2\nu) - \left( \frac{1}{2} q^{\text{R}} + L \right) \log(\pi) - \lambda(\boldsymbol{\xi})^\top (\boldsymbol{\xi}^2) \\ &\quad + \mathbf{1}^\top \psi(\boldsymbol{\xi}) + \left( \mathbf{y} - \frac{1}{2}\mathbf{1} \right)^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}| \\ &\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^{\text{G}} + p + m \right) \\ &\quad - \log(\mathcal{C}_{q^{\text{R}}, \nu + q^{\text{R}} - 1}) + \log(\mathcal{C}_{q^{\text{R}}, \nu + m + q^{\text{R}} - 1}) - \frac{1}{2} (\nu + m + q^{\text{R}} - 1) \log |B_{q(\boldsymbol{\Sigma}^{\text{R}})}| \\ &\quad + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^{\text{G}} + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^{\text{G}} + 1) \log(B_{q(\sigma_{u_\ell}^2)}) \\ &\quad - \sum_{r=1}^{q^{\text{R}}} \log(A_{\text{R}r}) + q^{\text{R}} \log \Gamma \left\{ \frac{1}{2} (\nu + q^{\text{R}}) \right\} \\ &\quad + \sum_{r=1}^{q^{\text{R}}} \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})})_{rr} \mu_{q(1/a_r^{\text{R}})} \\ &\quad - \frac{1}{2} (\nu + q^{\text{R}}) \sum_{r=1}^{q^{\text{R}}} \log(B_{q(a_r^{\text{R}})}) - \sum_{\ell=1}^L \log(A_{u_\ell}) \\ &\quad - \sum_{\ell=1}^L \log(B_{q(a_{u_\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u_\ell})} \mu_{q(1/\sigma_{u_\ell}^2)}, \end{aligned}$$

*Derivation:*

The logarithm of the lower bound on the marginal likelihood is given by

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \mathbf{a}^{\text{R}}, \mathbf{a}_u) \right. \\ &\quad \left. - \log q^*(\boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \mathbf{a}^{\text{R}}, \mathbf{a}_u) \right\}. \end{aligned}$$

According to the product restriction (2.21), our expression for the lower bound of the

marginal log-likelihood becomes:

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= E_q \{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) \} \\ &\quad + E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) \} \\ &\quad + E_q \{ \log p(\boldsymbol{\Sigma}^R | \mathbf{a}^R) - \log q^*(\boldsymbol{\Sigma}^R) \} E_q \{ p(\boldsymbol{\sigma}_u^2 | \mathbf{a}_u) - \log q^*(\boldsymbol{\sigma}_u^2) \} \\ &\quad + E_q \{ \log p(\mathbf{a}^R) - \log q^*(\mathbf{a}^R) \} + E_q \{ \log p(\mathbf{a}_u) - \log q^*(\mathbf{a}_u) \}. \end{aligned}$$

We are only interested in the first term (highlighted in darker colour) since it is the only difference between the Gaussian and Bernoulli response cases.

Replacing  $q^*$  by  $q^*(\cdot; \boldsymbol{\xi})$  we then have the following

$$\begin{aligned} \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) &= \mathbf{y}^\top \mathbf{C} \mathbf{v} - \mathbf{1}^\top \log(\mathbf{1} + \exp(\mathbf{C} \mathbf{v})) \\ &\geq \mathbf{y}^\top \mathbf{C} \mathbf{v} + \mathbf{v}^\top \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} \mathbf{v} - \frac{1}{2} \mathbf{1}^\top \mathbf{C} \mathbf{v} + \mathbf{1}^\top \psi(\boldsymbol{\xi}) \\ \log \underline{p}(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G; \boldsymbol{\xi}) &= (\mathbf{y} - \frac{1}{2} \mathbf{1})^\top \mathbf{C} \mathbf{v} + \mathbf{v}^\top \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} \mathbf{v} + \mathbf{1}^\top \psi(\boldsymbol{\xi}). \end{aligned}$$

Taking expectations,

$$\begin{aligned} E_{\underline{q}} \{ \log \underline{p}(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G; \boldsymbol{\xi}) \} &= (\mathbf{y} - \frac{1}{2} \mathbf{1})^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \text{tr} \left( \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \right. \\ &\quad \left. + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^\top \right) + \mathbf{1}^\top \psi(\boldsymbol{\xi}). \end{aligned}$$

## 2.D Optimal $q$ -densities derivation for Poisson semiparametric mixed models

Expressions for  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$

$$\begin{aligned} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} (\mathbf{D}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \mathbf{S})^\top, \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \left\{ -2 \text{vec}^{-1} \left( \left( \mathbf{D}_{\text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})} \mathbf{S} \right)^\top \right) \right\}^{-1}, \end{aligned}$$

where

$$\mathbf{S} \equiv E_q \{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) + \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2) \}$$

and the notation  $\mathbf{D}$  denotes derivative vector.

2.D. OPTIMAL  $q$ -DENSITIES DERIVATION FOR POISSON SEMIPARAMETRIC MIXED MODELS

---

*Derivation:*

First note that

$$\begin{aligned}\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) &= \mathbf{y}^\top (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top \log(\mathbf{y}!). \\ \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \sigma_u^2) &= \log \left[ (2\pi)^{-p/2} (\sigma_\beta^2 \mathbf{I}_p)^{-p/2} \exp\left(\frac{1}{2} \sigma_\beta^{-2} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right) \right. \\ &\quad \times (2\pi)^{-m/2} (\boldsymbol{\Sigma}^R)^{-m/2} \exp \left\{ -\frac{1}{2} (\mathbf{u}^R)^\top (\boldsymbol{\Sigma}^R)^{-1} \mathbf{u}^R \right\} \\ &\quad \left. \times \prod_{\ell=1}^L (2\pi)^{-q_\ell^G/2} (\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G})^{-q_\ell^G/2} \exp \left\{ \frac{1}{2} \sigma_{u\ell}^{-2} (\mathbf{u}_\ell^G)^\top \mathbf{u}_\ell^G \right\} \right].\end{aligned}$$

An explicit expression for  $S$  is

$$\begin{aligned}S &= \mathbf{y}^\top \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} - \mathbf{1}^\top \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal} \left( \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \right) \right\} \\ &\quad - \frac{1}{2} \text{tr} \left( \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^\top + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \right) - \frac{1}{2} \left( \sum_{i=1}^m n_i + \sum_{\ell=1}^L q_\ell^G + 1 \right) \log(2\pi) \\ &\quad - \frac{1}{2} \sum_{i=1}^m n_i E_q \{ \log(|\boldsymbol{\Sigma}^R|) \} - \frac{1}{2} \sum_{\ell=1}^L q_\ell^G E_q \{ \log(\sigma_{u\ell}^2) \} \\ &\quad - E_q \{ \log(\sigma_\beta^2) \} - \mathbf{1}^\top \log(\mathbf{y}!).\end{aligned}$$

Differentiate  $S$  with respect to  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  gives

$$\begin{aligned}d_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} S &= \mathbf{y}^\top \mathbf{C} d_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} - \mathbf{1}^\top \text{diag} \left\{ \exp(\mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) + \frac{1}{2} \text{diagonal} (\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top) \right\} \mathbf{C} d_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \\ &\quad - \frac{1}{2} \left( 2 \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^\top \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} d_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \right).\end{aligned}$$

Thus, by Theorem 6, Chapter 5 of Magnus and Neudecker (1995)

$$\begin{aligned}(\mathbf{D}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} S)^\top &= \mathbf{C}^\top \left[ \mathbf{y} - \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal} \left( \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \right) \right\} \right] \\ &\quad - \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}.\end{aligned}$$

Differentiate  $S$  with respect to  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  gives

$$\begin{aligned}d_{\text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})} S &= -\mathbf{1}^\top \text{diag} \left[ \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} Q(\mathbf{C}) \text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\} \right] \frac{1}{2} Q(\mathbf{C}) d_{\text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})} \\ &\quad - \frac{1}{2} \text{vec} \left( \mathbf{C}^\top \text{diag} \left[ \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal} \left( \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \right) \right\} \mathbf{C} \right. \right. \\ &\quad \left. \left. + \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \right] \right)^\top d_{\text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})},\end{aligned}$$

where  $Q(\mathbf{C}) = (\mathbf{C} \otimes \mathbf{1}^\top) \odot (\mathbf{1}^\top \otimes \mathbf{C})$  is defined in Theorem 2 of Section 4 of Wand (2014).

## 2.D. OPTIMAL $q$ -DENSITIES DERIVATION FOR POISSON SEMIPARAMETRIC MIXED MODELS

---

Therefore,

$$\text{vec}^{-1} \left\{ (\text{D}_{\text{vec}(\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})} S)^\top \right\} = \frac{1}{2} \left( \mathbf{C}^\top \text{diag} \left[ \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \mathbf{C}^\top) \right\} \right] \mathbf{C} + \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \right).$$

The remaining  $q^*$  densities involve straightforward adaptation of the derivations given in the Gaussian response case.

### 2.D.1 Derivation of the marginal log-likelihood lower bound

The expression for the marginal log-likelihood lower bound is

$$\begin{aligned} \log p(\mathbf{y}; q) &= \frac{1}{2} q^{\text{R}} (\nu + q^{\text{R}} - 1) \log(2\nu) - \left( \frac{1}{2} q^{\text{R}} + L \right) \log(\pi) \\ &\quad + \mathbf{y}^\top \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} - \mathbf{1}^\top \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \mathbf{C}^\top \right\} \\ &\quad - \mathbf{1}^\top \log(\mathbf{y}!) + \left( \mathbf{y} - \frac{1}{2} \mathbf{1} \right)^\top \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u}; \boldsymbol{\xi})} + \frac{1}{2} \log \|\boldsymbol{\Sigma}_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}\| \\ &\quad - \frac{1}{2} p \log(\sigma_\beta^2) - \frac{1}{2} \sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right) \\ &\quad - \log(\mathcal{C}_{q^{\text{R}}, \nu + q^{\text{R}} - 1}) + \log(\mathcal{C}_{q^{\text{R}}, \nu + m + q^{\text{R}} - 1}) \\ &\quad - \frac{1}{2} (\nu + m + q^{\text{R}} - 1) \log |B_{q(\boldsymbol{\Sigma}^{\text{R}})}| \\ &\quad + \sum_{\ell=1}^L \log \Gamma \left\{ \frac{1}{2} (q_\ell^{\text{G}} + 1) \right\} - \frac{1}{2} \sum_{\ell=1}^L (q_\ell^{\text{G}} + 1) \log(B_{q(\sigma_{u\ell}^2)}) \\ &\quad + q^{\text{R}} \log \Gamma \left\{ \frac{1}{2} (\nu + q^{\text{R}}) \right\} + \sum_{r=1}^{q^{\text{R}}} \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})})_{rr} \mu_{q(1/a_r^{\text{R}})} \\ &\quad - \frac{1}{2} (\nu + q^{\text{R}}) \sum_{r=1}^{q^{\text{R}}} \log(B_{q(a_r^{\text{R}})}) - \sum_{\ell=1}^L \log(A_{u\ell}) \\ &\quad - \sum_{\ell=1}^L \log(B_{q(a_{u\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u\ell})} \mu_{q(1/\sigma_{u\ell}^2)} \\ &\quad + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^{\text{G}} + p + m \right) - \sum_{r=1}^{q^{\text{R}}} \log(A_{Rr}), \end{aligned}$$

*Derivation:*

The derivation is very similar to that of the Bernoulli case.

## Chapter 3

# Streamlining Mean Field

# Variational Bayes Algorithms And Three-Level Model Extensions

*Technology is so much fun but we can drown in our technology. The fog of information can drive out knowledge.*

Daniel J. Boorstin

### 3.1 Introduction

Previously we introduced the concept of MFVB algorithms that aim to approximate the intractable posterior by a factorised distribution which can be represented by a directed acyclic graph, and optimisation of the factorised approximating posterior can be decomposed into local computations that involve only neighbouring nodes. Recent examples of these include Algorithm 3 of Ormerod and Wand (2010) and Algorithms 3 and 5 of Luts *et al.* (2014) which use naïve matrix inversion for the effects covariance matrix updates to support MFVB fitting of longitudinal and multilevel data. Following the naïve implementation by these authors, we presented the naïve version of MFVB algorithms which are not fully computationally optimised because of their cubic dependence on the number of groups.

---

The main content of this chapter is published as: Lee, C. Y. Y. and Wand, M. P. (2015). Streamlined mean field variational Bayes for longitudinal and multilevel data analysis. *Biometrical Journal article*. DOI:10.1002/bimj.201500007. This research has been accepted for presentation at four conferences.

In this chapter, we develop MFVB algorithms for fitting and inference in arbitrarily large longitudinal and multilevel models that are streamlined in terms of number of operations and storage. The number of operations is linear in the number of groups at each level, which constitutes a major improvement over the naïve implementations. The essence of our approach is to take advantage of the sparseness of the matrix requiring inversion and streamline its inversion. This involves omitting correlations of the estimated effects between groups, which are rarely of interest. It is likely that our resultant streamlined algorithms provide the fastest ever means of approximate Bayesian analyses of very large longitudinal and multilevel datasets.

Recently, Tan and Nott (2013) introduce a different strategy for improving the efficiency of variational methods for longitudinal and multilevel models. The product restriction of Tan and Nott (2013) is such that the random effects for each group appear in a separate factor and their *partially noncentered parameterisation* strategy is aimed at improved accuracy in the face of such a restriction. We do not impose their restriction in our variational approximation for the random group effects. Fitting and inference proceeds efficiently via a streamlined approach that takes advantage of the inherent blocked structure of the effects covariance matrix. Our algorithms correspond to the algorithms of Tan and Nott (2013) in which their tuning parameter, that captures correlations between fixed and random effects, does not have to be estimated and, instead, is specified optimally. Our streamlined algorithms are simpler due to minimal product restrictions in the MFVB approximation and apply to a richer class of models.

The next section presents the computational obstacles involved in the naïve MFVB algorithms. Section 3.4 gives details of our novel streamlined approach to sparse covariance estimation via matrix permutation and decomposition, and is the centrepiece of this chapter. In Section 3.5 we provide numerical evidence of the efficiency gain for using the streamlined approach over the naïve approach. Illustration of real data examples is presented in Section 3.6. Section 3.7 extends the streamlined variational algorithms to cater for higher-level Bayesian semiparametric mixed models and concluding remarks are given in Section 3.8.

## 3.2 Computational challenges in naïve MFVB algorithms

Recall from Chapter 2 that the general compact form of a two-level Bayesian generalised semiparametric mixed model is

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \mathbf{u} &\sim \text{Dist}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \mathbf{R}), \quad \boldsymbol{\beta} \sim N(0, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}_p), \\ \mathbf{u}|\sigma_{u_1}^2, \dots, \sigma_{u_L}^2, \boldsymbol{\Sigma}^R &\sim N\left(\mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u_\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma}^R \end{bmatrix}\right), \end{aligned} \tag{3.1}$$



$$\begin{aligned}
 \Sigma^{\mathbf{R}} | a_1^{\mathbf{R}}, \dots, a_{q^{\mathbf{R}}}^{\mathbf{R}} &\sim \text{Inverse-Wishart} \left( \nu + q^{\mathbf{R}} - 1, 2\nu \text{diag}(1/a_1^{\mathbf{R}}, \dots, 1/a_{q^{\mathbf{R}}}^{\mathbf{R}}) \right), \\
 a_r^{\mathbf{R}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{\mathbf{R}r}^2 \right), \quad r = 1, \dots, q^{\mathbf{R}}, \\
 \sigma_{u\ell}^2 | a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/a_{u\ell} \right), \quad a_{u\ell} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{u\ell}^2 \right).
 \end{aligned}$$

The notation ‘‘Dist’’ denotes the type of response distribution including Gaussian, Student- $t$ , Bernoulli or Poisson and  $\mathbf{R}$  represents the additional parameters corresponding to the errors, e.g.  $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}$  for the Gaussian response mixed model. We illustrated in the previous chapter approximate Bayesian fitting of model (3.1) via iterative coordinate descent MFVB algorithms. While perhaps not intuitively obvious, these algorithms are challenged by the increasing volume and complexity of longitudinal and multilevel data as we now elaborate.

A computational problem arising consistently across Algorithms 3, 4 and 6 lies in the update expression for the effects covariance matrix. Take the Gaussian response model as an example,

$$\Sigma_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \left( \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^T \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^{\mathbf{G}}})_{1 \leq \ell \leq L} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes M_{q((\Sigma^{\mathbf{R}})^{-1})} \end{bmatrix} \right)^{-1}. \quad (3.2)$$

This update expression requires storage and inversion of a large sparse matrix of dimension  $(p + q^{\mathbf{R}} m + \sum_{\ell=1}^L q_\ell^{\mathbf{G}}) \times (p + q^{\mathbf{R}} m + \sum_{\ell=1}^L q_\ell^{\mathbf{G}})$ , where  $p$  is the number of predictors,  $q_\ell^{\mathbf{G}}$  ( $1 \leq \ell \leq L$ ) is the size of the  $\ell$ th penalised spline, typically in the range 15 – 40 regardless sample size, and  $m$  is the number of groups that can be arbitrarily large. For the real data examples we consider in this chapter,  $m$  does not pose serious problems. However, other studies involve a much larger  $m$ , the *Six Cities Study of Air Pollution and Health* described in Fitzmaurice *et al.* (2012) with  $m = 13\,379$  is a practical example. Thus, without doubt the number of groups  $m$  dominates the dimension of  $\Sigma_{q(\boldsymbol{\beta}, \mathbf{u})}$  in many practical situations. Henceforth, if  $p = 10$ ,  $q^{\mathbf{R}} = 2$ ,  $m = 10\,000$ ,  $\ell = 1$  and  $q_1^{\mathbf{G}} = 25$ , then the dimension of the matrix requiring storage and inversion would exceed  $20\,000 \times 20\,000$ . In addition, it is well known from numerical linear algebra that *naïve* computation of  $\Sigma_{q(\boldsymbol{\beta}, \mathbf{u})}$  is  $O(m^3)$ , that is cubic dependence on the number of groups. Hence, direct implementation of Algorithms 3, 4 and 6 can be very costly, or even infeasible, in large longitudinal and multilevel studies. Our aim in the next few sections is to *streamline* these MFVB algorithms by removing the computational obstacles involving update expressions such as (3.2) for most practical values of  $m$  and, thus, maximising the benefits of MFVB methods for approximate Bayesian inference.

### 3.3 Predictor structure and matrix notation

Here we introduce new matrix notation for our predictor structure that will benefit our presentation of the streamlined approach later on. Our predictor structure corresponds to the set up notation of Section 2 of Zhao *et al.* (2006) with the spatial correlation structure omitted. We first partition  $\beta$ ,  $\mathbf{X}$ ,  $\mathbf{u}$  and  $\mathbf{Z}$  into random group effects (superscript R) and general spline components (superscript G) as follows:

$$\beta \equiv \begin{bmatrix} \beta^{\text{R}} \\ \beta^{\text{G}} \end{bmatrix}, \quad \mathbf{X} \equiv [\mathbf{X}^{\text{R}} \ \mathbf{X}^{\text{G}}], \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}^{\text{R}} \\ \mathbf{u}^{\text{G}} \end{bmatrix} \quad \text{and} \quad \mathbf{Z} \equiv [\mathbf{Z}^{\text{R}} \ \mathbf{Z}^{\text{G}}],$$

which then leads to

$$\mathbf{X}\beta + \mathbf{Z}\mathbf{u} = \mathbf{X}^{\text{R}}\beta^{\text{R}} + \mathbf{X}^{\text{G}}\beta^{\text{G}} + \mathbf{Z}^{\text{R}}\mathbf{u}^{\text{R}} + \mathbf{Z}^{\text{G}}\mathbf{u}^{\text{G}}. \quad (3.3)$$

The random group effects design matrices are of the form

$$\mathbf{X}^{\text{R}} \equiv \begin{bmatrix} \mathbf{X}_1^{\text{R}} \\ \vdots \\ \mathbf{X}_m^{\text{R}} \end{bmatrix} \quad \text{and} \quad \mathbf{Z}^{\text{R}} \equiv \text{blockdiag}(\mathbf{X}_i^{\text{R}})_{1 \leq i \leq m},$$

where the  $\mathbf{X}_i^{\text{R}}$ ,  $1 \leq i \leq m$ , are  $n_i \times q^{\text{R}}$  design matrices. The random group effects vector is such that

$$\mathbf{u}^{\text{R}} = \begin{bmatrix} \mathbf{u}_1^{\text{R}} \\ \vdots \\ \mathbf{u}_m^{\text{R}} \end{bmatrix} \quad \text{and} \quad \text{Cov}(\mathbf{u}^{\text{R}}) = \mathbf{I}_m \otimes \Sigma^{\text{R}},$$

where  $\Sigma^{\text{R}}$  is an unstructured  $q^{\text{R}} \times q^{\text{R}}$  covariance matrix. Thus,  $\mathbf{X}^{\text{R}}\beta^{\text{R}} + \mathbf{Z}^{\text{R}}\mathbf{u}^{\text{R}}$  corresponds to random intercepts and slopes for repeated measures data on  $m$  groups with sample sizes  $n_1, \dots, n_m$ .

The matrices  $\mathbf{X}^{\text{G}}$  and  $\mathbf{Z}^{\text{G}}$  are general design matrices corresponding to the fixed effects vector  $\beta^{\text{G}}$  and spline coefficients vector  $\mathbf{u}^{\text{G}}$  respectively. Typically,  $\mathbf{X}^{\text{G}}$  contains polynomial functions of a continuous predictor that enter the model as a penalised spline. The  $\mathbf{Z}^{\text{G}}$  matrix would then contain spline basis functions of the same predictor. Mixed model based penalised spline fitting of (3.3) (e.g. Ruppert *et al.*, 2003) involves modelling a smooth, but otherwise unspecified, function  $f$  according to:

$$f(x) = \beta_x x + \sum_{k=1}^{q^{\text{G}}} u_k^{\text{G}} z_k(x), \quad u_k^{\text{G}} \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2),$$

where  $\{z_k : 1 \leq k \leq q^G\}$  are spline bases of size  $q^G$  and  $\sigma_u^2$  is the penalised parameter for the spline coefficients  $u_1^G, \dots, u_{q^G}^G$ . A common choice for the  $z_k$  is the suitably linearly transformed cubic O’Sullivan spline, as described in Section 4 of Wand and Ormerod (2008). The  $\mathbf{Z}^G \mathbf{u}^G$  term can be further decomposed according to

$$\mathbf{Z}^G \mathbf{u}^G \equiv \sum_{\ell=1}^L \mathbf{Z}_\ell^G \mathbf{u}_\ell^G \quad \text{with} \quad \text{Cov}(\mathbf{u}^G) = \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}). \quad (3.4)$$

Here the  $\mathbf{u}_\ell^G$  are  $q_\ell^G \times 1$  random spline coefficient vectors for  $1 \leq \ell \leq L$ , so the blocks of  $\text{Cov}(\mathbf{u}^G)$  correspond to the decomposition of  $\mathbf{u}^G$ . The expression in (3.4) is in keeping with spline penalisation in multi-predictor semiparametric regression models (e.g. Ruppert *et al.*, 2003). Further, it is useful to define  $\mathbf{C}^G \equiv [\mathbf{X} \ \mathbf{Z}^G]$  and partition  $\mathbf{y}$ ,  $\mathbf{X}$ ,  $\mathbf{Z}$  and  $\mathbf{C}^G$  row-wise corresponding to the groups in  $\mathbf{X}^R$ :

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_m \end{bmatrix} \quad \text{and} \quad \mathbf{C}^G = \begin{bmatrix} \mathbf{C}_1^G \\ \vdots \\ \mathbf{C}_m^G \end{bmatrix}.$$

Here  $\mathbf{y}_i \equiv [y_{i1} \dots y_{in_i}]^\top$  denotes the  $n_i \times 1$  vector of responses for the  $i$ th group. The matrices  $\mathbf{X}_i$ ,  $\mathbf{Z}_i$  and  $\mathbf{C}_i^G$  are defined in the same manner.

### 3.4 Streamlining mean field variational Bayes algorithms

The direct approach to obtaining the optimal  $q$ -density moments for all key parameters in Algorithms 3, 4 and 6 involves computation of  $\Sigma_{q(\beta, \mathbf{u})}$ . As mentioned in Section 3.2, the matrix  $\Sigma_{q(\beta, \mathbf{u})}$  increases in dimension as the number of groups  $m$  increases, and the time complexity is as high as  $O(m^3)$ , making variational algorithms become computationally infeasible in large longitudinal and multilevel studies. We get around this by exploiting the fact that most of the matrix requiring inversion is block-diagonal. Hence we propose a streamlined approach based around block decomposition of a matrix, using the Gaussian response model (2.11) as a demonstrating example. First we re-express (3.2) in a simpler

matrix form as given below:

$$\begin{aligned} \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} + & \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes M_{q((\Sigma^R)^{-1})} \end{bmatrix} \\ = & \begin{bmatrix} \mathbf{F} & \mathbf{G}_1 & \mathbf{G}_2 & \cdots & \mathbf{G}_m \\ \mathbf{G}_1^\top & \mathbf{H}_1^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{G}_2^\top & \mathbf{0} & \mathbf{H}_2^{-1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{G}_m^\top & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{H}_m^{-1} \end{bmatrix}, \end{aligned} \quad (3.5)$$

where the submatrices  $\mathbf{F}$ ,  $\mathbf{G}_i$  and  $\mathbf{H}_i$  are defined as follows:

$$\begin{aligned} \mathbf{F} & \equiv \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}^G)^\top \mathbf{C}^G + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G}) \end{bmatrix}, \\ \mathbf{G}_i & \equiv \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}_i^G)^\top \mathbf{X}_i^R, \\ \text{and } \mathbf{H}_i & \equiv \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{X}_i^R)^\top \mathbf{X}_i^R + M_{q((\Sigma^R)^{-1})} \right\}^{-1}. \end{aligned}$$

The block-diagonal structure in the bottom right submatrix of (3.5), the contribution from the random group effects, is crucial as it enables us to reduce the inversion to  $O(m)$  operations.

**Result 3.4.1.** Let  $\mathbf{A}$  be a  $m \times m$  matrix,  $\mathbf{B}$  a  $m \times n$  matrix,  $\mathbf{C}$  an  $n \times m$  matrix and  $\mathbf{D}$  an  $n \times n$  matrix. Suppose that the partitioned matrix  $\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$  is nonsingular, its inverse can then be obtained by solving a system of simultaneous linear equations in matrix form

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{A}} & \tilde{\mathbf{B}} \\ \tilde{\mathbf{C}} & \tilde{\mathbf{D}} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix},$$

where the submatrices are defined as follows:

$$\begin{aligned} \tilde{\mathbf{A}} & = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} \\ \tilde{\mathbf{B}} & = \tilde{\mathbf{A}}\mathbf{B}\mathbf{D}^{-1}, \quad \tilde{\mathbf{C}} = -\mathbf{D}^{-1}\mathbf{C}\tilde{\mathbf{A}} \\ \tilde{\mathbf{D}} & = \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\tilde{\mathbf{A}}\mathbf{B}\mathbf{D}^{-1}, \end{aligned}$$

assuming the inverse matrices  $\mathbf{D}^{-1}$  and  $(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})$  above exist.

The standard Result 3.4.1 (Harville, 2008) on inversion of a block-partitioned matrix leads to the inverse of (3.5) equalling

$$\Sigma_{q(\beta, \mathbf{u})} \equiv \begin{bmatrix} \Sigma_{q(\beta, \mathbf{u}^G)} & \Lambda_{q(\beta, \mathbf{u}^G, \mathbf{u}_1^R, \dots, \mathbf{u}_m^R)} \\ \Lambda_{q(\beta, \mathbf{u}^G, \mathbf{u}_1^R, \dots, \mathbf{u}_m^R)}^\top & \Sigma_{q(\mathbf{u}^R)} \end{bmatrix}, \quad (3.6)$$

where

$$\Sigma_{q(\beta, \mathbf{u}^G)} \equiv \left( \mathbf{F} - \sum_{i=1}^m \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \right)^{-1}$$

and  $\Lambda_{q(\beta, \mathbf{u}^G, \mathbf{u}_1^R, \dots, \mathbf{u}_m^R)} \equiv \left[ -\Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_1 \mathbf{H}_1 \quad \dots \quad -\Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_m \mathbf{H}_m \right].$

The dimension of the matrix  $\Sigma_{q(\beta, \mathbf{u}^G)}$  is  $(p + \sum_{\ell=1}^L q_\ell^G) \times (p + \sum_{\ell=1}^L q_\ell^G)$  and is relatively straightforward to compute. The summation term renders this whole process as  $O(m)$ . This matrix is all we require to obtain variability bands for the penalised regression splines.

To find the variability estimates to accompany group-specific mean estimates, we need  $\Sigma_{q(\mathbf{u}^R)}$  which is not a block-diagonal matrix. However, since the covariance between the fitted values of two different groups is rarely of interest, it suffices to compute and store the diagonal blocks

$$\Sigma_{q(\mathbf{u}_i^R)} \equiv \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i, \quad 1 \leq i \leq m. \quad (3.7)$$

Since  $\Sigma_{q(\beta, \mathbf{u}^G)}$ ,  $\mathbf{G}_i$  and  $\mathbf{H}_i$  have dimensions much smaller than  $m$ , the complexity of the matrix calculations required in these submatrices does not increase as  $m$  increases. Therefore, the calculations with the highest order of complexity are the calculations of the  $m$  relevant submatrices of  $\Lambda_{q(\beta, \mathbf{u}^G, \mathbf{u}_i^R)}$  and  $\Sigma_{q(\mathbf{u}_i^R)}$ . This renders the whole process as  $O(m)$ , representing an improvement of order  $m^2$  over the naïve approach to matrix inversion. As the number of groups increases, the improvements due to streamlining become enormous. It is worth mentioning that automatic sparse methods could be used, without requiring detailed algebra, but the downside of this is a possibly sub-optimal implementation compared to knowing the block-diagonal structure.

### 3.4.1 Streamlining update expressions involving $\Sigma_{q(\beta, \mathbf{u})}$

We see in Algorithm 3 that the update expressions for the mean vector of the coefficient estimates and the scale parameter of the error variance involve the matrix  $\Sigma_{q(\beta, \mathbf{u})}$ . Henceforth we now work towards a streamlined alternative to these updates.

Recall from Algorithm 3 that the naïve update for  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  is

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{y}.$$

By replacing  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  with (3.6) and partitioning  $\mathbf{C}$  into  $[(\mathbf{C}^\mathbf{G})^\top (\mathbf{Z}^\mathbf{R})^\top]^\top$ , we can partition the above update into two expressions corresponding to  $(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})$  and  $\mathbf{u}^\mathbf{R}$  respectively:

$$\begin{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} \\ \boldsymbol{\mu}_{q(\mathbf{u}^\mathbf{R})} \end{bmatrix} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \begin{bmatrix} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} & \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G}, \mathbf{u}_1^\mathbf{R}, \dots, \mathbf{u}_m^\mathbf{R})} \\ \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G}, \mathbf{u}_1^\mathbf{R}, \dots, \mathbf{u}_m^\mathbf{R})}^\top & \boldsymbol{\Sigma}_{q(\mathbf{u}^\mathbf{R})} \end{bmatrix} \begin{bmatrix} (\mathbf{C}^\mathbf{G})^\top \mathbf{y} \\ (\mathbf{Z}^\mathbf{R})^\top \mathbf{y} \end{bmatrix}.$$

Specifically, after some simple algebraic manipulations, we obtain the streamlined update for  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})}$  and  $\boldsymbol{\mu}_{q(\mathbf{u}^\mathbf{R})}$  respectively:

$$\begin{aligned} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \left( \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} (\mathbf{C}^\mathbf{G})^\top \mathbf{y} + \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G}, \mathbf{u}_1^\mathbf{R}, \dots, \mathbf{u}_m^\mathbf{R})} (\mathbf{Z}^\mathbf{R})^\top \mathbf{y} \right) \quad (3.8) \\ &= \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} \left( (\mathbf{C}^\mathbf{G})^\top \mathbf{y} - \sum_{i=1}^m \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^\mathbf{R})^\top \mathbf{y}_i \right) \\ \text{and } \boldsymbol{\mu}_{q(\mathbf{u}^\mathbf{R})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \left( \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G}, \mathbf{u}_1^\mathbf{R}, \dots, \mathbf{u}_m^\mathbf{R})}^\top (\mathbf{C}^\mathbf{G})^\top \mathbf{y} + \boldsymbol{\Sigma}_{q(\mathbf{u}^\mathbf{R})} (\mathbf{Z}^\mathbf{R})^\top \mathbf{y} \right). \end{aligned}$$

Using (3.6) and (3.7), we can further break down  $\boldsymbol{\mu}_{q(\mathbf{u}^\mathbf{R})}$  into components corresponding to the  $i$ th group only:

$$\begin{aligned} \boldsymbol{\mu}_{q(\mathbf{u}_i^\mathbf{R})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \left\{ -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^\mathbf{R})^\top \mathbf{y}_i \right. \quad (3.9) \\ &\quad \left. + \left( \mathbf{H}_i (\mathbf{X}_i^\mathbf{R})^\top \mathbf{y}_i + \mathbf{H}_i \mathbf{G}_i^\top \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} \sum_{i=1}^m \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^\mathbf{R})^\top \mathbf{y}_i \right) \right\} \\ &= \mathbf{H}_i \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{X}_i^\mathbf{R})^\top \mathbf{y}_i - \mathbf{G}_i^\top \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^\mathbf{G})} \right\}. \end{aligned}$$

Similarly, recall from Algorithm 3 that the naïve update for  $B_{q(\sigma_\varepsilon^2)}$  is

$$B_{q(\sigma_\varepsilon^2)} \leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) + \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\}.$$

By replacing  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  with (3.8) and partitioning  $\mathbf{C}$  into  $[(\mathbf{C}^\mathbf{G})^\top (\mathbf{Z}^\mathbf{R})^\top]^\top$ , we can rewrite

$\mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)})$  in the following equivalent streamlined form:

$$\begin{aligned} \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \mathbf{C}^G \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} + \begin{bmatrix} \mathbf{X}_1^R \boldsymbol{\mu}_{q(\mathbf{u}_1^R)} \\ \vdots \\ \mathbf{X}_m^R \boldsymbol{\mu}_{q(\mathbf{u}_m^R)} \end{bmatrix} \quad \text{and} \quad (3.10) \\ \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}) &= \text{tr}\left\{(\mathbf{C}^G)^\top \mathbf{C}^G \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}\right\} - 2\mu_{q(1/\sigma_\varepsilon^2)}^{-1} \sum_{i=1}^m \text{tr}\left(\mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}\right) \\ &\quad + \sum_{i=1}^m \text{tr}\left\{(\mathbf{X}_i^R)^\top \mathbf{X}_i^R \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}\right\}, \end{aligned}$$

and the streamlined update for  $B_{q(\sigma_\varepsilon^2)}$  then follows.

### 3.4.2 Streamlining lower bound expression involving $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$

The marginal log-likelihood lower bound  $\log p(\mathbf{y}; q)$  for Algorithm 3 is obtained via straightforward, albeit long-winded, algebra. The only change in the lower bound expression for the streamlined approach is to compute  $\log|\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|$  more efficiently, which we now justify. It depends on the following result concerning determinant of block-diagonal matrices:

$$\begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = |\mathbf{D}| |\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}|, \quad (3.11)$$

where  $\mathbf{D}$  is a square matrix and is invertible (e.g. Theorem 13.3.8 of Harville, 2008). From (3.11) we then get

$$\begin{aligned} \log|\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}| &= -\log \begin{vmatrix} \mathbf{F} & \mathbf{G}_1 & \mathbf{G}_2 & \cdots & \mathbf{G}_m \\ \mathbf{G}_1^\top & \mathbf{H}_1^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{G}_2^\top & \mathbf{0} & \mathbf{H}_2^{-1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{G}_m^\top & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{H}_m^{-1} \end{vmatrix} \\ &= -\sum_{i=1}^m \log|\mathbf{H}_i^{-1}| - \log|\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}^{-1}|. \end{aligned}$$

We illustrate step-by-step transformations of update expressions in Algorithm 3 concerning the  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  term for the Gaussian response model. It is, however, relatively straightforward to extrapolate the steps discussed and apply to other response models that were covered in Chapter 2. Hence we omit the details and directly present the streamlined algorithms for much faster and memory-efficient MFVB approximate fitting

of longitudinal and multilevel models. These are presented as Algorithms 7, 8 and 9.

Algorithms 7, 8 and 9 are an alternative to (restricted) maximum likelihood and best prediction-based fitting of large longitudinal and multilevel models, the latter being employed by popular software such as PROC MIXED in SAS (SAS Institute Inc., 2013) and `lmer()` in the R package `lme4` (Pinheiro *et al.*, 2014). Leaving aside the fact that the latter are based on frequentist inference paradigms, whilst Algorithms 7, 8 and 9 are intrinsically Bayesian, it is worth pointing out that these algorithms are purely matrix algebraic which has advantages such as stability, speed, parallelisability and online fitting. On the other hand, (restricted) maximum likelihood estimation of covariance matrices involves multi-dimensional, nonlinear optimisation (e.g. Wolfinger *et al.*, 1994). Such optimisation procedures become computationally burdensome for very large longitudinal and multilevel datasets, as exemplified by the computing times given in Table 3.5 as part of the speed comparisons given in Section 3.5. Earlier covariance estimation procedures, such as the minimum norm quadratic unbiased estimator method of Rao (1972) and the moment-based approach described in Section 3 of Goldstein (1986), are faster but have drawbacks such as positive definiteness not being guaranteed.

### 3.5 Computational speed: naïve versus streamlined

We conducted a simulation study to assess the computational speed of the naïve Algorithm 3 against that of the streamlined Algorithm 7. We generated 100 datasets according to the following simulation setting, focusing on the Gaussian response model.

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R)x_{ij} + f(s_{ij}), \sigma_\varepsilon^2),$$

where

$$f(s) = 1 - \frac{13}{5\sqrt{2\pi}} e^{-(s-0.15)^2/0.2} - (2.3s - 0.07s^2) + 0.5 \{1 - \Phi(s; 0.8, 0.07)\}$$

and  $s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1)$ ,  $x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1)$  and  $[u_{0i}^R \ u_{1i}^R]^\top | \Sigma^R \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \Sigma^R)$ , for  $1 \leq i \leq m$ ,  $1 \leq j \leq n_i$ . The true parameter values are:

$$\beta_0 = 0.58, \quad \beta_x = 1.89, \quad \sigma_\varepsilon^2 = 0.04 \quad \text{and} \quad \Sigma^R = \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix}.$$

The number of groups is varied,  $m \in \{100, 500, 2500, 12500\}$ , with the within-group sample sizes  $n_i$  ranged between 10 and 20.

Both the naïve Algorithm 3 and streamlined Algorithm 7 were implemented in Fortran



**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mu_{q(1/\sigma_{u_\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\mu_{q(1/b_i)} > 0$ ,  $1 \leq i \leq N$  and  $\mu_{q(\nu)} > 0$ .

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega)} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L}(\mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G}) \end{bmatrix}$$

$$N \leftarrow \sum_{i=1}^m n_i \quad ; \quad \text{Gaussian: } \Psi = \mathbf{I} \quad ; \quad \text{Student-t: } \Psi \leftarrow \text{diag}_{1 \leq i \leq N}(\mu_{q(1/b_i)})$$

$$\mathbf{S} \leftarrow \mathbf{0} \quad ; \quad \mathbf{s} \leftarrow \mathbf{0}$$

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

$$\begin{aligned} \mathbf{G}_i &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}_i^G)^\top \Psi_i \mathbf{X}_i^R \\ \mathbf{H}_i &\leftarrow \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{X}_i^R)^\top \Psi_i \mathbf{X}_i^R + \mathbf{M}_{q((\Sigma^R)^{-1})} \right\}^{-1} \\ \mathbf{S} &\leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \quad ; \quad \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^R)^\top \mathbf{y}_i \end{aligned}$$

**1. Update multivariate normal  $q^*(\beta, \mathbf{u}^G)$  parameters:**

$$\begin{aligned} \Sigma_{q(\beta, \mathbf{u}^G)} &\leftarrow \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}^G)^\top \Psi \mathbf{C}^G + \mathbf{M}_{q(\Omega)} - \mathbf{S} \right\}^{-1} \\ \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta, \mathbf{u}^G)} \left\{ (\mathbf{C}^G)^\top \mathbf{y} - \mathbf{s} \right\} \end{aligned}$$

**Update multivariate normal  $q^*(\mathbf{u}_i^R)$  parameters:**

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

$$\begin{aligned} \Sigma_{q(\mathbf{u}_i^R)} &\leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i \\ \boldsymbol{\mu}_{q(\mathbf{u}_i^R)} &\leftarrow \mathbf{H}_i \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{X}_i^R)^\top \mathbf{y}_i - \mathbf{G}_i^\top \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} \right\} \end{aligned}$$

**2. Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  and inverse-Gamma  $q^*(a_\varepsilon)$  parameters:**

$$\mathbf{R} \leftarrow \mathbf{y} - \mathbf{C}^G \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} - \begin{bmatrix} \mathbf{X}_1^R \boldsymbol{\mu}_{q(\mathbf{u}_1^R)} \\ \vdots \\ \mathbf{X}_m^R \boldsymbol{\mu}_{q(\mathbf{u}_m^R)} \end{bmatrix}$$

$$\begin{aligned} B_{q(\sigma_\varepsilon^2)} &\leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left[ \mathbf{R}^\top \Psi \mathbf{R} + \text{tr}\{(\mathbf{C}^G)^\top \Psi \mathbf{C}^G \Sigma_{q(\beta, \mathbf{u}^G)}\} \right. \\ &\quad \left. + \sum_{i=1}^m \text{tr}\{(\mathbf{X}_i^R)^\top \Psi_i \mathbf{X}_i^R \Sigma_{q(\mathbf{u}_i^R)}\} \right. \\ &\quad \left. - 2(\mu_{q(1/\sigma_\varepsilon^2)})^{-1} \sum_{i=1}^m \text{tr}(\mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^G)}) \right] \end{aligned}$$

$$\mu_{q(1/\sigma_\varepsilon^2)} \leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2})$$

**3. Update inverse-Gamma  $q^*(b_i)$  and  $q^*(\nu)$  parameters for the Student- $t$  response model:**

$$\zeta \leftarrow \text{diagonal}\{\mathbf{C}^G \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} (\mathbf{C}^G)^\top\}$$

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

$$\begin{aligned} \zeta_i &\leftarrow \zeta_i - 2 \text{diagonal}\{\mathbf{X}_i^R (\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i)^\top (\mathbf{C}_i^G)^\top\} \\ \zeta_i &\leftarrow \zeta_i + \text{diagonal}\{\mathbf{X}_i^R \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)} (\mathbf{X}_i^R)^\top\} \end{aligned}$$

For  $1 \leq i \leq N$ :

$$\begin{aligned} B_{q(b_i)} &\leftarrow \frac{1}{2} \left[ \mu_{q(\nu)} + \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})_i^2 + \zeta_i \right\} \right] \\ \mu_{q(1/b_i)} &\leftarrow \frac{1}{2} (\mu_{q(\nu)} + 1) / B_{q(b_i)} \\ \mu_{q(\log b_i)} &\leftarrow \log(B_{q(b_i)}) - \text{digamma} \left\{ \frac{1}{2} (\mu_{q(\nu)} + 1) \right\} \end{aligned}$$

$$C_1 \leftarrow \sum_{i=1}^N (\mu_{q(\log(b_i))} + \mu_{q(1/b_i)})$$

$$\mu_{q(\nu)} \leftarrow \exp \{ \log \mathcal{F}(1, N, C_1, \nu_{\min}, \nu_{\max}) - \log \mathcal{F}(0, N, C_1, \nu_{\min}, \nu_{\max}) \}$$

**4. Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\begin{aligned} \mu_{q(1/a_{u\ell})} &\leftarrow 1 / (\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2}) \\ \mu_{q(\sigma_{u\ell}^2)} &\leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)}) + 2 \mu_{q(1/a_{u\ell})}} \end{aligned}$$

**5. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\boldsymbol{\Sigma}^R)$  parameters:**

For  $r = 1, \dots, q^R$ :

$$\begin{aligned} B_{q(a_r^R)} &\leftarrow \nu \{ \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \}_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a_r^R)} \leftarrow \frac{1}{2} (\nu + q^R) / B_{q(a_r^R)} \\ \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} &\leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}) + 2 \nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)}) \\ \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} &\leftarrow (\nu + m + q^R - 1) \mathbf{B}_{q(\boldsymbol{\Sigma}^R)}^{-1} \end{aligned}$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

**Update  $q$ -density covariance matrix for the  $i$ th group:**

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

$$\boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^G, \mathbf{u}_i^R)} \equiv E_q \left[ \left( \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix} - \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \right) (\mathbf{u}_i^R - \boldsymbol{\mu}_{q(\mathbf{u}_i^R)})^\top \right] \leftarrow -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i$$

---

Algorithm 7: Streamlined mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed models with Gaussian and Student- $t$  responses.

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix and  $\boldsymbol{\xi}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega)} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} \end{bmatrix}$$

$$\mathbf{S} \leftarrow \mathbf{0} \ ; \ \mathbf{s} \leftarrow \mathbf{0}$$

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

$$\begin{aligned} \mathbf{G}_i &\leftarrow 2(\mathbf{C}_i^G)^\top \text{diag}\{\lambda(\boldsymbol{\xi}_i)\} \mathbf{X}_i^R \\ \mathbf{H}_i &\leftarrow \left\{ 2(\mathbf{X}_i^R)^\top \text{diag}\{\lambda(\boldsymbol{\xi}_i)\} \mathbf{X}_i^R + \mathbf{M}_{q((\Sigma^R)^{-1})} \right\}^{-1} \\ \mathbf{S} &\leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \ ; \ \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^R)^\top (\mathbf{y}_i - \tfrac{1}{2} \mathbf{1}) \end{aligned}$$

**1. Update Multivariate Normal  $q^*(\boldsymbol{\beta}, \mathbf{u}^G)$  parameters:**

$$\begin{aligned} \Sigma_{q(\boldsymbol{\beta}, \mathbf{u}^G)} &\leftarrow \left\{ 2(\mathbf{C}^G)^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C}^G + \mathbf{M}_{q(\Omega)} - \mathbf{S} \right\}^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} &\leftarrow 2 \Sigma_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \left\{ (\mathbf{C}^G)^\top (\mathbf{y} - \tfrac{1}{2} \mathbf{1}) - \mathbf{s} \right\} \end{aligned}$$

**Update Multivariate Normal  $q^*(\mathbf{u}^R)$  parameters:**

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

$$\begin{aligned} \Sigma_{q(\mathbf{u}_i^R)} &\leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i \\ \boldsymbol{\mu}_{q(\mathbf{u}_i^R)} &\leftarrow \mathbf{H}_i \left\{ 2(\mathbf{X}_i^R)^\top (\mathbf{y}_i - \tfrac{1}{2} \mathbf{1}) - \mathbf{G}_i^\top \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \right\} \\ \boldsymbol{\xi}^2 &\leftarrow \text{diagonal}\{ \mathbf{C}^G (\Sigma_{q(\boldsymbol{\beta}, \mathbf{u}^G)} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)}^\top) (\mathbf{C}^G)^\top \} \end{aligned}$$

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

$$\begin{aligned} \boldsymbol{\xi}_i^2 &\leftarrow \boldsymbol{\xi}_i^2 + 2 \text{diagonal}\{ \mathbf{C}^G (-\Sigma_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top) (\mathbf{X}_i^R)^\top \} \\ \boldsymbol{\xi}_i^2 &\leftarrow \boldsymbol{\xi}_i^2 + \text{diagonal}\{ (\mathbf{X}_i^R)^\top (\Sigma_{q(\mathbf{u}_i^R)} + \boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top) (\mathbf{X}_i^R)^\top \} \end{aligned}$$

**2. Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\begin{aligned} \mu_{q(1/a_{u\ell})} &\leftarrow 1 / (\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2}) \\ \mu_{q(\sigma_{u\ell}^2)} &\leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}_\ell^G)}) + 2 \mu_{q(1/a_{u\ell})}} \end{aligned}$$

**3. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\Sigma^R)$  parameters:**

 For  $\ell = 1, \dots, q^R$ :

$$B_{q(a^R)} \leftarrow \nu(\mathbf{M}_{q((\Sigma^R)^{-1})})_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a^R)} \leftarrow \frac{1}{2}(\nu + q^R)/B_{q(a^R)}$$

$$\mathbf{B}_{q(\Sigma^R)} \leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}) + 2\nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_q^R)})$$

$$\mathbf{M}_{q((\Sigma^R)^{-1})} \leftarrow (\nu + \sum_{i=1}^m n_i + q^R - 1) \mathbf{B}_{q(\Sigma^R)}^{-1}$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

**Update  $q$ -density covariance matrix for the  $i$ th group:**

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

$$\boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^G, \mathbf{u}_i^R)} \equiv E_q \left[ \left( \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix} - \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \right) (\mathbf{u}_i^R - \boldsymbol{\mu}_{q(\mathbf{u}_i^R)})^\top \right] \leftarrow -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i$$


---

Algorithm 8: Streamlined mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Bernoulli response.

77. We also compared Algorithm 7 with contemporary software for fitting the frequentist version of (3.12). This was achieved using the function `gamm()` in the R package `mgcv` (Wood, 2014). Note that `gamm()` uses the function `lme()` in the R package `nlme` (Pinheiro *et al.*, 2014) as a fitting engine. All computations were performed on a Mac OS X laptop with a 2.6 GHz Intel Core i5 processor and 8 GBytes of random access memory.

Table 3.5 summarises the average (standard error) computing times over 100 runs and shows the practical benefits of the streamlined MFVB approach. As  $m$  increases the average computing time for the naïve approach increases rapidly from 0.2 seconds to almost 25 minutes, compared with about a quarter to half of a second for the streamlined approach. For  $m = 12500$ , the naïve approach failed due to required storage of  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$  exceeding memory restrictions for typical 2015 personal computing environments. However, the streamlined approach just took over 1.8 seconds on average to compute. The impressive speed gains in the streamlined approach are clearly reflected in the ratios of naïve over streamlined and the ratios of `gamm()` over streamlined with respect to the average computing times. In situations where a dataset has a large number of groups, the streamlined approach is more than 3500 times faster than the naïve approach and 1700 times faster than the `gamm()` approach.

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\mathbf{w}_{q(\beta, \mathbf{u})}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries and  $\boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)}$ , a  $(p + \sum_{\ell=1}^L q_\ell^G) \times 1$  vector.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega)} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G})_{1 \leq \ell \leq L} \end{bmatrix}$$

$$\mathbf{S} \leftarrow \mathbf{0} \ ; \ \mathbf{s} \leftarrow \mathbf{0}$$

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

$$\begin{aligned} \mathbf{G}_i &\leftarrow (\mathbf{C}_i^G)^\top \text{diag}\{(\mathbf{w}_{q(\beta, \mathbf{u})})_i\} \mathbf{X}_i^R \\ \mathbf{H}_i &\leftarrow \left\{ (\mathbf{X}_i^R)^\top \text{diag}\{(\mathbf{w}_{q(\beta, \mathbf{u})})_i\} \mathbf{X}_i^R + \mathbf{M}_{q((\Sigma^R)^{-1})} \right\}^{-1} \\ \mathbf{S} &\leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \ ; \ \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{X}_i^R)^\top \{\mathbf{y}_i - (\mathbf{w}_{q(\beta, \mathbf{u})})_i\} \end{aligned}$$

**1. Update Multivariate Normal  $q^*(\beta, \mathbf{u}^G)$  parameters:**

$$\begin{aligned} \Sigma_{q(\beta, \mathbf{u}^G)} &\leftarrow \{(\mathbf{C}^G)^\top \boldsymbol{\Psi} \mathbf{C}^G + \mathbf{M}_{q(\Omega)} - \mathbf{S}\}^{-1} \\ \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} &\leftarrow \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} + \Sigma_{q(\beta, \mathbf{u}^G)} \left\{ (\mathbf{C}^G)^\top \mathbf{y} - \mathbf{s} - \text{diagonal}(\mathbf{M}_{q(\Omega)}) \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} \right\} \\ \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)}^{\text{old}} &\leftarrow \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} \end{aligned}$$

**Update Multivariate Normal  $q^*(\mathbf{u}_i^R)$  parameters:**

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

$$\begin{aligned} \Sigma_{q(\mathbf{u}_i^R)} &\leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i \\ \boldsymbol{\mu}_{q(\mathbf{u}_i^R)} &\leftarrow \boldsymbol{\mu}_{q(\mathbf{u}_i^R)} + \mathbf{H}_i \left[ 2(\mathbf{X}_i^R)^\top \{\mathbf{y}_i - (\mathbf{w}_{q(\beta, \mathbf{u})})_i\} - \mathbf{G}_i^\top (\boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} - \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)}^{\text{old}}) \right] \end{aligned}$$

$$\boldsymbol{\zeta} \leftarrow \text{diagonal}\{\mathbf{C}^G \Sigma_{q(\beta, \mathbf{u}^G)} (\mathbf{C}^G)^\top\}$$

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

$$\begin{aligned} \boldsymbol{\zeta}_i &\leftarrow \boldsymbol{\zeta}_i - 2 \text{diagonal}\{\mathbf{X}_i^R (\Sigma_{q(\beta, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i)^\top (\mathbf{C}_i^G)^\top\} \\ \boldsymbol{\zeta}_i &\leftarrow \boldsymbol{\zeta}_i + \text{diagonal}\{\mathbf{X}_i^R \Sigma_{q(\mathbf{u}_i^R)} (\mathbf{X}_i^R)^\top\} \end{aligned}$$

$$\mathbf{w}_{q(\beta, \mathbf{u})} \leftarrow \exp \left( \mathbf{C}^G \boldsymbol{\mu}_{q(\beta, \mathbf{u}^G)} + \begin{bmatrix} \mathbf{X}_1^R \boldsymbol{\mu}_{q(\mathbf{u}_1^R)} \\ \vdots \\ \mathbf{X}_m^R \boldsymbol{\mu}_{q(\mathbf{u}_m^R)} \end{bmatrix} + \frac{1}{2} \boldsymbol{\zeta} \right)$$

**4. Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:**

For  $\ell = 1, \dots, L$ :

$$\begin{aligned}\mu_{q(1/a_{u\ell})} &\leftarrow 1/(\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2}) \\ \mu_{q(\sigma_{u\ell}^2)} &\leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell^G)}) + 2\mu_{q(1/a_{u\ell})}}\end{aligned}$$

**5. Update inverse-Gamma  $q^*(a_r^R)$  and inverse-Wishart  $q^*(\boldsymbol{\Sigma}^R)$  parameters:**

For  $r = 1, \dots, q^R$ :

$$\begin{aligned}B_{q(a_r^R)} &\leftarrow \nu(\mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})})_{rr} + A_{Rr}^{-2} \quad ; \quad \mu_{q(a_r^R)} \leftarrow \frac{1}{2}(\nu + q^R)/B_{q(a_r^R)} \\ \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} &\leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^R)} \boldsymbol{\mu}_{q(\mathbf{u}_i^R)}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^R)}) + 2\nu \text{diag}(\mu_{q(1/a_1^R)}, \dots, \mu_{q(1/a_{q^R}^R)}) \\ \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} &\leftarrow (\nu + m + q^R - 1) \mathbf{B}_{q(\boldsymbol{\Sigma}^R)}^{-1}\end{aligned}$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

**Update  $q$ -density covariance matrix for the  $i$ th group:**

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

$$\boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}^G, \mathbf{u}_i^R)} \equiv E_q \left[ \left( \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^G \end{bmatrix} - \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \right) (\mathbf{u}_i^R - \boldsymbol{\mu}_{q(\mathbf{u}_i^R)})^\top \right] \leftarrow -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^G)} \mathbf{G}_i \mathbf{H}_i$$

---

Algorithm 9: Streamlined mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Poisson response.

### 3.6. REAL DATA APPLICATIONS

$m$	Naïve	Streamlined	<code>gamm</code>	Naïve/Stream.	<code>gamm</code> /Stream.
100	0.198 (0.023)	0.027 (0.003)	0.17 (0.010)	7.4	5.8
500	14.482 (1.704)	0.106 (0.021)	7.075 (0.165)	136.8	76.2
2500	1488.757 (108.311)	0.419 (0.056)	739.687 (17.793)	3550.4	1765.4
12500	Failed	1.766 (0.096)	Failed	N/A	N/A

Table 3.1: Average (standard error) elapsed of the computing times in seconds for the simulation described in the text, using the naïve Algorithm 3, streamlined Algorithm 7 and `gamm()` function in R `mgcv` package. The ratio of naïve over streamlined and the ratio of `gamm()` over streamlined are also presented.

## 3.6 Real data applications

We now provide illustration of our streamlined MFVB methodology through a series of re-analysis of real data examples. Throughout this section, we standardise the response and continuous variables to have zero means and unit standard deviations. We use a normal prior for  $\beta$  and a half-Cauchy prior for  $\sigma_R^2$ ,  $\sigma_\varepsilon^2$  and  $\sigma_u^2$ , with the values of the hyperparameters all set to 10 000.

### 3.6.1 Application to smoking data

The first example is based on the perinatal health dataset from the United States National Center for Health Statistics, which is discussed and analysed by Abrevaya (2006). A 10% random sub-sample is provided by Rabe-Hesketh and Skrondal (2008) and our analysis is of this subset. The data have a two-level structure with 8604 births as units at Level 1 and 3978 mothers as groups at Level 2. There are an average of 2.2 births per mother. The following variables are given in Table 3.2.

In this example, study variables can either vary at the birth level (therefore also at the mother level) or at the mother level only. For instance, `smoke` is a Level-1 variable as maternal smoking status can change from one pregnancy to the next, whereas `black` is a Level-2 variable as mother’s ethnicity remains unchanged between pregnancies.

Motivated by a report from the United States Surgeon General: “Infants born to women who smoke during pregnancy have a lower average birthweight and are more likely to be small for gestational age than infants born to women who do not smoke ...” Abrevaya (2006) examined the effect of smoking on birthweight for women in the United States between 1990 and 1998 using a matched panel data approach. Here we re-analyse their data using a different approach via streamlined MFVB. The main study factor is maternal smoking status and the response of interest is infant’s birthweight. The following Bayesian Gaussian semiparametric mixed model with a random intercept for each mother was fitted

### 3.6. REAL DATA APPLICATIONS

Variable	Description
momid	mother identifier
birwt	birthweight in grams
gestat	infant's gestational age in weeks
mage	mother's age at the birth of the infant in years
smoke	indicator for mother smoking during pregnancy
male	indicator for infant being male
married	indicator for mother being married
hsgrad	indicator for mother having some college education, but not degree
somecoll	indicator for mother having graduated from college
black	indicator for mother being black
kessner2	indicator for Kessner index equalling 2
kessner3	indicator for Kessner index equalling 3
novisit	indicator for no prenatal care visit
pretri2	indicator for first prenatal care visit having occurred in second trimester
pretri3	indicator for first prenatal care visit having occurred in third trimester

Table 3.2: Description of the United States National Center for Health Statistics perinatal health data as presented in Abrevaya (2006).

via Algorithm 7:

$$\begin{aligned}
 \text{birwt}_{ij} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\stackrel{\text{ind.}}{\sim} N \left( \beta_0 + u_i^R + \beta_1 \text{smoke}_{ij} + \beta_2 \text{mage}_{ij} + \beta_3 \text{male}_{ij} + \beta_4 \text{married}_{ij} \right. \\
 &\quad + \beta_5 \text{hsgrad}_{ij} + \beta_6 \text{somecoll}_{ij} + \beta_7 \text{collgrad}_{ij} + \beta_8 \text{black}_{ij} \\
 &\quad + \beta_9 \text{kessner2}_{ij} + \beta_{10} \text{kessner3}_{ij} + \beta_{11} \text{novisit}_{ij} \\
 &\quad \left. + \beta_{12} \text{pretri2}_{ij} + \beta_{13} \text{pretri3}_{ij} + f(\text{gestat}_{ij}), \sigma_\varepsilon^2 \right), \\
 \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p), \quad u_i^R | \sigma_R^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_R^2), \quad \sigma_R^2 \sim \text{Half-Cauchy}(0, A_R), \\
 \sigma_\varepsilon^2 &\sim \text{Half-Cauchy}(0, A_\varepsilon), \quad u_k^G | \sigma_u^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2), \quad \sigma_u^2 \sim \text{Half-Cauchy}(0, A_u) \quad (3.12)
 \end{aligned}$$

where

$$f(\text{gestat}) = \beta_{14} \text{gestat} + \sum_{k=1}^{q^G} u_k^G z_k(\text{gestat})$$

is a penalised spline function for gestational age. Here  $z_1(\cdot), \dots, z_{q^G}(\cdot)$  is a set of O'Sullivan spline basis functions and  $\sigma_u^2$  represents the amount of penalisation of the spline coefficients  $u_1^G, \dots, u_{q^G}^G$  as described in Section 1.3.

Figures 3.1 and 3.2 allow a visual assessment of the MFVB-based approximate posterior density functions against a MCMC benchmark for (3.12). MCMC was fitted using Stan with a burn-in of size 5000, a post burn-in of size 5000 with a thinning factor of 5. The



accuracy of approximate posterior density functions of  $\beta$  and  $f(Q_k)$  is excellent, ranging from 95% to 98%, while it is about 75% to 82% for the variance parameters  $\sigma_u^2$  and  $\sigma_\varepsilon^2$  respectively. Inspection of Figure 3.1 shows that the MFVB fits and pointwise 95% credible sets are very close to those obtained using MCMC. The MCMC fits took 4.1 days, while the MFVB fits took 1.4 seconds. This represents more than a two hundred thousand-fold improvement in computing time.

### 3.6.2 Application to student assessment data

The second example is based on a dataset from the 2000 Program for International Student Assessment conducted by the Organisation for Economic Cooperation and Development. The survey assessed education attainment of 15 years old students in 43 counties, with an emphasis on reading proficiency. Our analysis is of the United States sample of the full dataset, provided by Rabe-Hesketh and Skrondal (2008). The data have a two-level structure with 2069 observations as units at Level 1 and 148 schools as groups at Level 2. The following variables are given in Table 3.3.

Variable	Description
<code>idschool</code>	school identifier
<code>passread</code>	indicator for being proficient in reading
<code>isei</code>	International Socioeconomic Index (SES)
<code>college</code>	indicator for highest education level by either parent being college
<code>oneforeign</code>	indicator for one parent being foreign born
<code>twoforeign</code>	indicator for both parents being foreign born
<code>language</code>	indicator for test language (English) being spoken at home

Table 3.3: Description of the 2000 Program for International Student Assessment conducted by the Organisation for Economic Cooperation and Development.

We treat reading proficiency as the response variable, an indicator taking the value 1 to indicate proficient and 0 otherwise. The effect of socio-economic status has considerable interest in education, we therefore model its effect the most flexible way via a penalised spline. The following Bayesian semiparametric random intercept logistic regression model for student  $i$  in school  $j$  was fitted via Algorithm 8:

$$\begin{aligned}
 \text{passread}_{ij} | \beta, \mathbf{u} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left( \text{logit}^{-1} \left\{ \beta_0 + u_i^R + \beta_1 \text{education}_{ij} + \beta_2 \text{oneforeign}_{ij} \right. \right. \\
 &\quad \left. \left. + \beta_3 \text{bothforeign}_{ij} + \beta_4 \text{language}_{ij} + f(\text{isei}_{ij}) \right\} \right), \\
 \beta &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p), \quad u_i^R | \sigma_R^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_R^2), \quad \sigma_R^2 \sim \text{Half-Cauchy}(0, A_R), \\
 u_k^G | \sigma_u^2 &\stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2) \quad \text{and} \quad \sigma_u^2 \sim \text{Half-Cauchy}(0, A_u).
 \end{aligned} \tag{3.13}$$

The function  $f$  is a penalised spline for the International Socio-economic Index, defined

analogously in (3.12). Figures 3.1 and 3.3 allow a visual assessment of the MFVB-based approximate posterior density functions against a MCMC benchmark for (3.13). MCMC was fitted using Stan with a burn-in of size 5000, a post burn-in of size 5000 with a thinning factor of 5. The accuracy of the approximate posterior density functions of  $\beta$  and  $f(Q_k)$  is good, ranging from 87% to 94%. In contrast, the MFVB approximation is poor for the variance parameter  $\sigma_R^2$ , as we previously observed in the Bernoulli simulation study results. Nonetheless, inspection of Figure 3.1 shows that the MFVB fits and pointwise 95% credible sets are close to those obtained using MCMC. The MCMC fits took 57 minutes, while the MFVB fits took 1.2 minutes.

### 3.6.3 Application to German health care data

The final example is based on the German Socio-Economic Panel dataset, an ongoing annual household survey that was started in 1984. The data are presented in Winkelmann (2004), and a subset of the original data is provided by Rabe-Hesketh and Skrondal (2008) and our analysis is of this subset. The data have a two-level structure with 32 837 observations as units at Level 1 and 9197 patients as groups at Level 2. Each patient has one to five observations. The data have a few variables relating to the usage of health service, and of particular interest is the number of visits to a doctor during the previous three months. The following variables are given in Table 3.4.

Variable	Description
<code>id</code>	patient identifier
<code>numvisit</code>	number of visits to a doctor during the 3 months before the interview
<code>male</code>	indicator for being a male
<code>married</code>	indicator for being married
<code>sport</code>	indicator for being actively in sports
<code>goodhealth</code>	indicator for being in good health
<code>badhealth</code>	indicator for being in bad health
<code>welfare</code>	indicator for one receiving welfare payments
<code>fulltime</code>	indicator for being employed full-time
<code>unemploy</code>	indicator for being unemployed
<code>winter</code>	indicator for being interviewed in winter quarter
<code>spring</code>	indicator for being interviewed in spring quarter
<code>fall</code>	indicator for being interviewed in fall quarter

Table 3.4: Description of the German health care data as presented in Winkelmann (2004).

We treat the number of doctor visits as the response variable, and of particular interest is the age effects which we modelled flexibly via a penalised spline. The following Bayesian semiparametric random intercept Poisson regression model for patient  $i$  in observation  $j$  was fitted via Algorithm 9:

$$\begin{aligned}
 \text{numvisit}_{ij} | \boldsymbol{\beta}, \mathbf{u} &\stackrel{\text{ind.}}{\sim} \text{Poisson} \left( \exp \left\{ \beta_0 + u_i^{\text{R}} + \beta_1 \text{male}_{ij} + \beta_2 \text{married}_{ij} + \beta_3 \text{sport}_{ij} \right. \right. \\
 &\quad + \beta_4 \text{goodhealth}_{ij} + \beta_5 \text{badhealth}_{ij} + \beta_6 \text{welfare}_{ij} \\
 &\quad + \beta_7 \text{fulltime}_{ij} + \beta_8 \text{unemploy}_{ij} + \beta_9 \text{winter}_{ij} \\
 &\quad \left. \left. + \beta_{10} \text{spring}_{ij} + \beta_{11} \text{fall}_{ij} + f(\text{isei}_{ij}) \right\} \right), \\
 \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p), \quad u_i^{\text{R}} | \sigma_{\text{R}}^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_{\text{R}}^2), \quad \sigma_{\text{R}}^2 \sim \text{Half-Cauchy}(0, A_{\text{R}}), \\
 u_k^{\text{G}} | \sigma_u^2 &\stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2) \quad \text{and} \quad \sigma_u^2 \sim \text{Half-Cauchy}(0, A_u). \tag{3.14}
 \end{aligned}$$

The function  $f$  is a penalised spline for age, defined analogously in (3.12). Figures 3.1 and 3.4 allow a visual assessment of the MFVB-based approximate posterior density functions against a MCMC benchmark for (3.13). MCMC was fitted using Stan with a burn-in of size 5000, a post burn-in of size 5000 with a thinning factor of 5. The accuracy of the approximate posterior density functions for all parameters is excellent, ranging from 80% to 97%. Inspection of Figure 3.4 shows that both the MFVB and MCMC spline fits are virtually identical. The MCMC fits took 1.6 hours while the MFVB fits took 48 seconds.

### 3.7 Extension to three-level semiparametric mixed models

A great advantage of MFVB algorithms is *modularity*. By this we means that concepts like main effects, interaction effects, higher-order random effects and spline regression can be viewed as modules that can be put together into an almost endless variety of statistical models. All that is required is relatively straightforward modifications on the structures of the general design matrices and effects covariance matrix in order to accommodate larger and more complicated mixed models. Consider a general three-level mixed models with the Gaussian response, and with simplicity in mind, we confine the model to the case of a single predictor, say  $x$ ,

$$\begin{aligned}
 y_{ijk} &= \beta_0 + \beta_x x_{ijk} + \{(u_{0ij}^{\text{RL2}} + u_{0i}^{\text{RL3}}) + (u_{1ij}^{\text{RL2}} + u_{1i}^{\text{RL3}}) x_{ijk}\} + \varepsilon_{ijk}, \tag{3.15} \\
 [u_{0ij}^{\text{RL2}} \quad u_{1ij}^{\text{RL2}}]^\top &\stackrel{\text{ind.}}{\sim} N \left( \mathbf{0}, \boldsymbol{\Sigma}^{\text{RL2}} = \begin{bmatrix} (\sigma_{u_0}^{\text{RL2}})^2 & \sigma_{u_0, u_1}^{\text{RL2}} \\ \sigma_{u_0, u_1}^{\text{RL2}} & (\sigma_{u_1}^{\text{RL2}})^2 \end{bmatrix} \right), \\
 [u_{0i}^{\text{RL3}} \quad u_{1i}^{\text{RL3}}]^\top &\stackrel{\text{ind.}}{\sim} N \left( \mathbf{0}, \boldsymbol{\Sigma}^{\text{RL3}} = \begin{bmatrix} (\sigma_{u_0}^{\text{RL3}})^2 & \sigma_{u_0, u_1}^{\text{RL3}} \\ \sigma_{u_0, u_1}^{\text{RL3}} & (\sigma_{u_1}^{\text{RL3}})^2 \end{bmatrix} \right), \\
 \text{and} \quad \varepsilon_{ijk} &\stackrel{\text{ind.}}{\sim} N(0, \sigma_\varepsilon^2), \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij},
 \end{aligned}$$

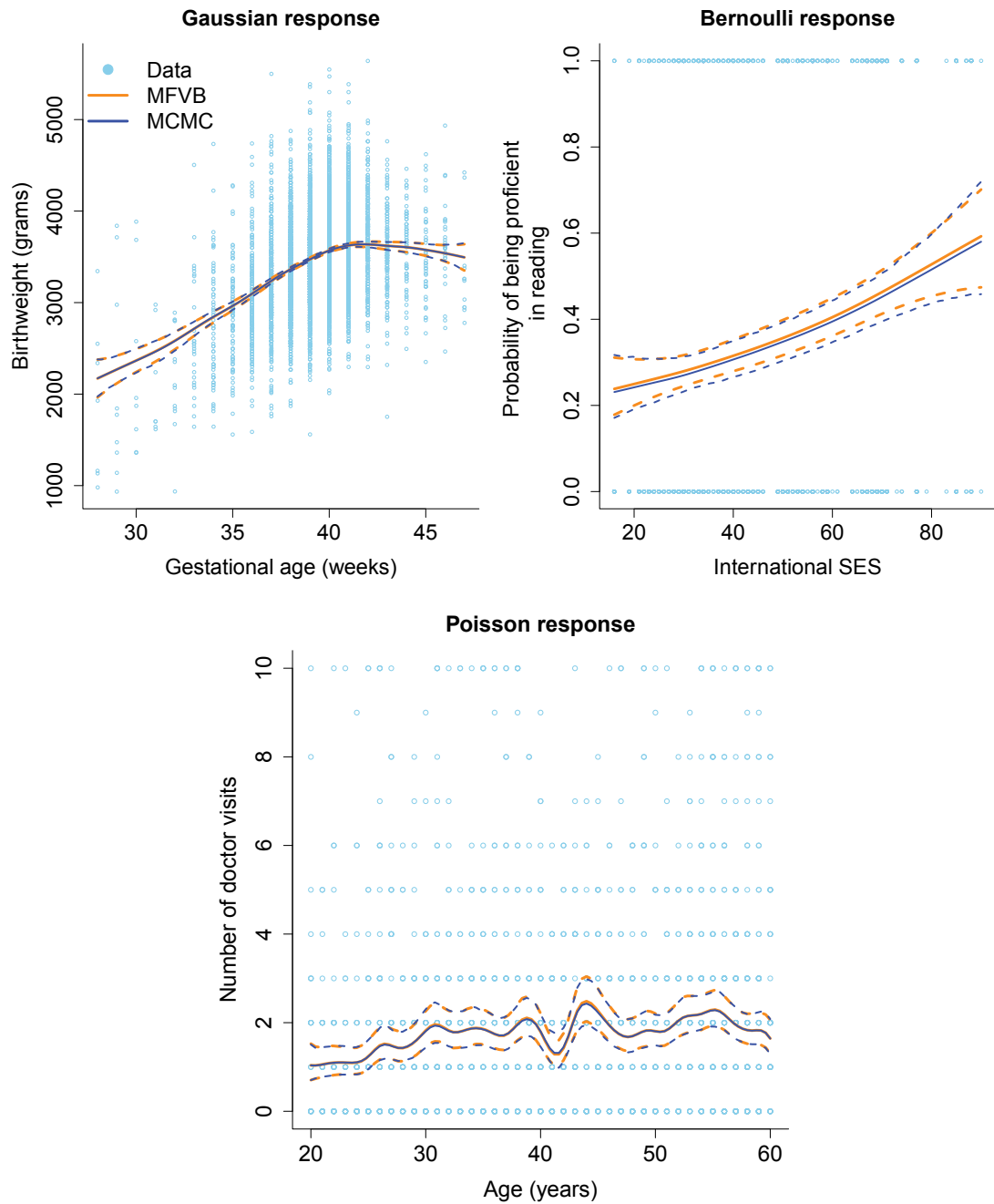


Figure 3.1: Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the two-level Bayesian semiparametric mixed models to the real datasets. The sky blue circles represent the real data.

### 3.7. EXTENSION TO THREE-LEVEL SEMIPARAMETRIC MIXED MODELS

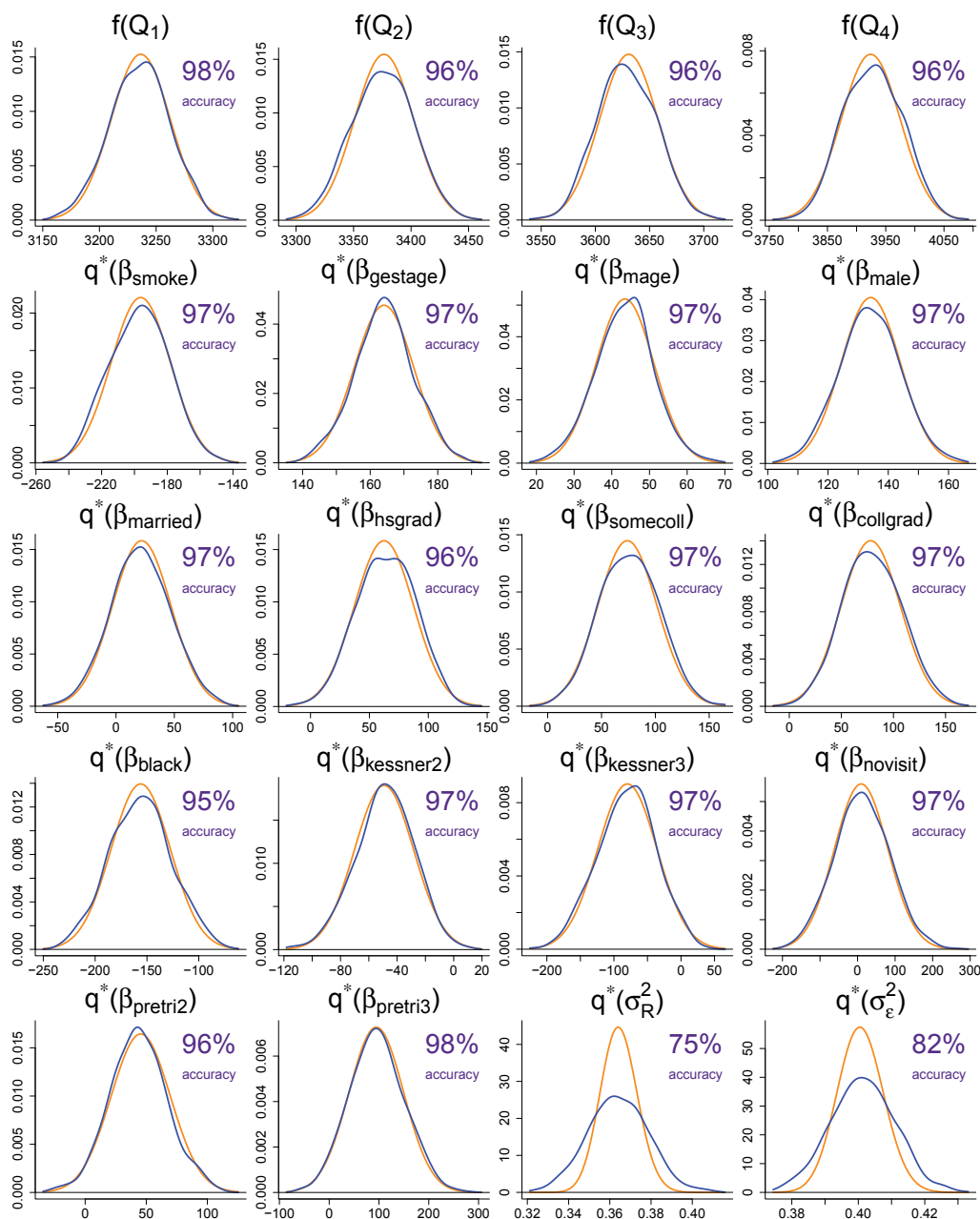


Figure 3.2: Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Gaussian response model to the smoking data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark.

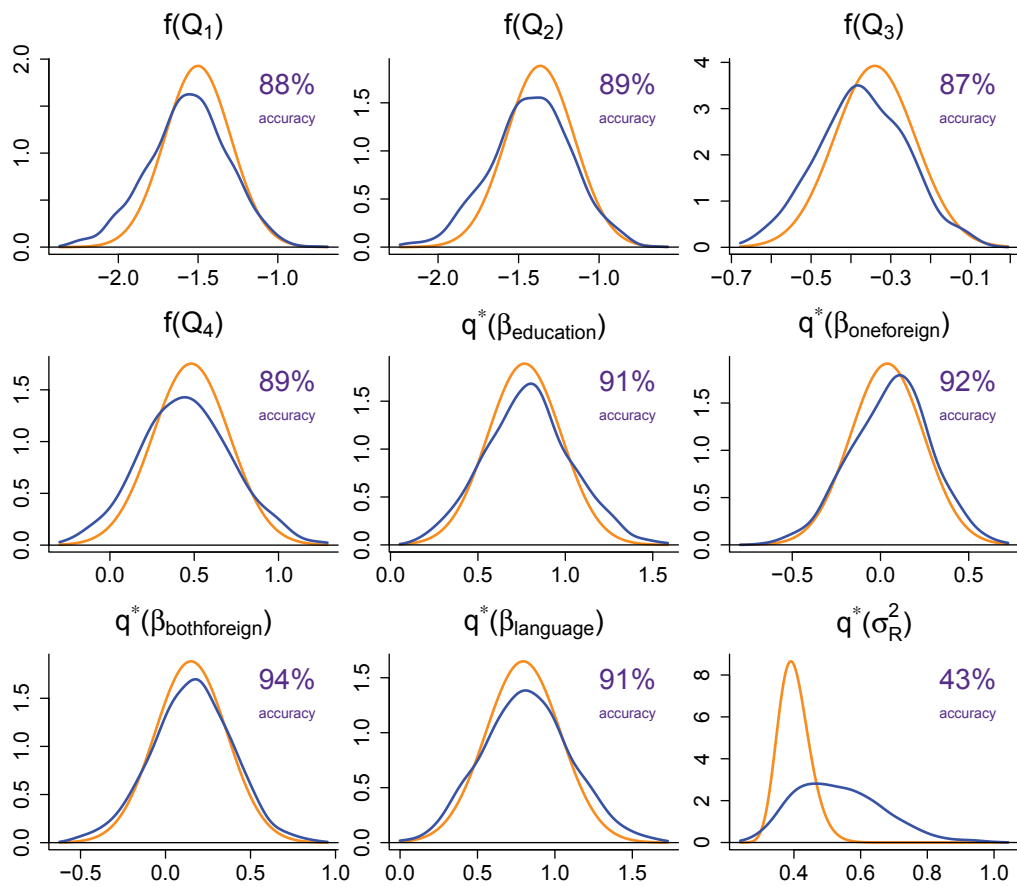


Figure 3.3: Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Bernoulli response model to the Program of International Student Assessment data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark.

### 3.7. EXTENSION TO THREE-LEVEL SEMIPARAMETRIC MIXED MODELS

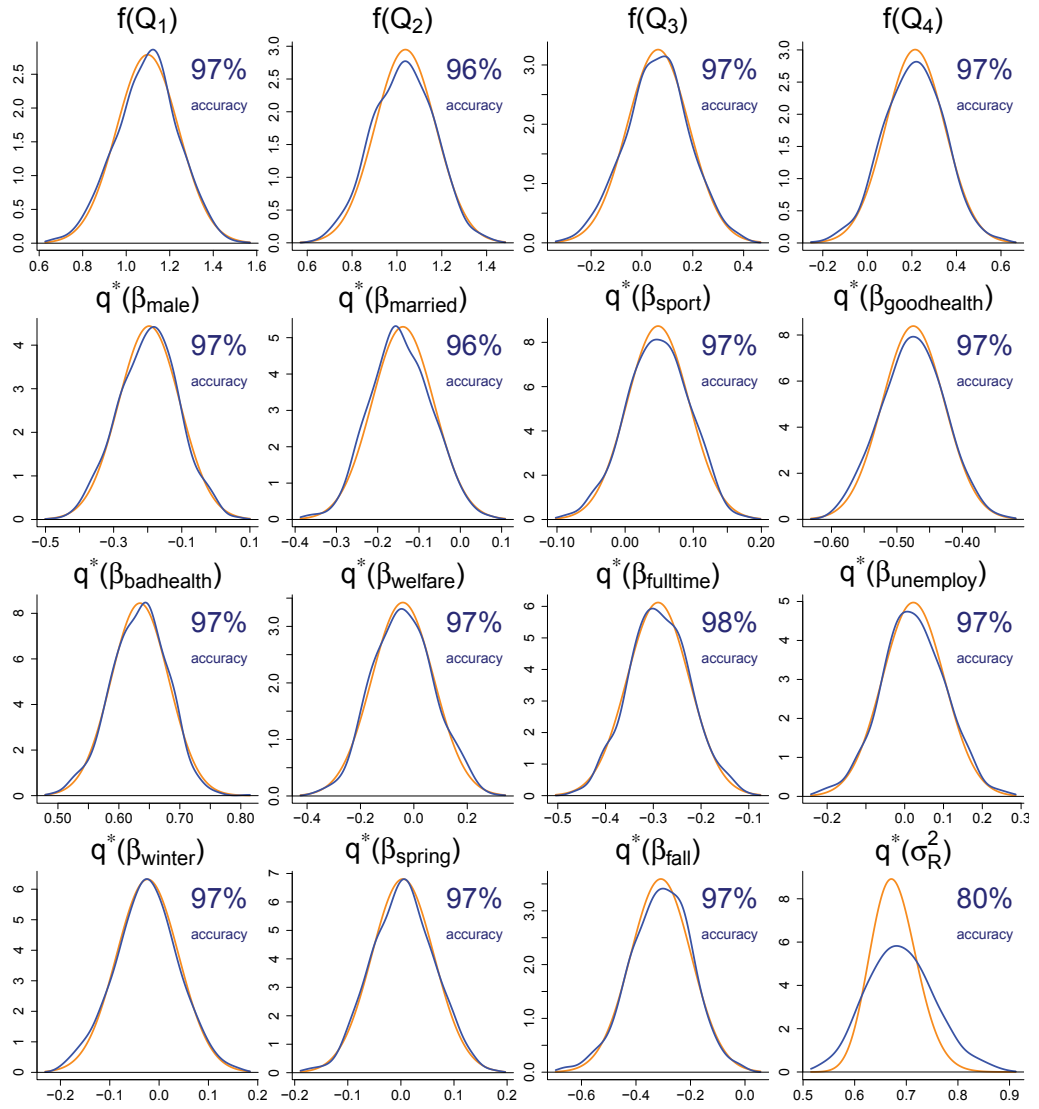


Figure 3.4: Approximate posterior density functions for the model parameters obtained via MFVB approximation (orange curves) and MCMC (blue curves) for the Poisson response model to the German health care data. The accuracy score on the top right of each plot represents the accuracy of MFVB approximation compared against the MCMC benchmark.

### 3.7. EXTENSION TO THREE-LEVEL SEMIPARAMETRIC MIXED MODELS

---

where  $i, j, k$  denote the *nested* indexing for the group level (Level 3, superscript RL3), subgroup level (Level 2, superscript RL2) and unit level (Level 1) respectively. The superscript notation is deliberately chosen to be in accordance with Section 3.3. The term  $y_{ijk}$  and  $x_{ijk}$  denote the response and predictor variables for the  $k$ th individual nested within the  $j$ th subgroup and the  $i$ th group,  $\beta_0$  and  $\beta_x$  are the so-called fixed effects,  $u_{0ij}^{\text{RL2}}$  and  $u_{1ij}^{\text{RL2}}$  are the so-called random intercepts and random slopes at the subgroup level, and similarly,  $u_{0i}^{\text{RL3}}$  and  $u_{1i}^{\text{RL3}}$  are the random intercepts and random slopes at the group level. Here we use the random effects to model the dependence of  $y$ s within groups and subgroups. It is reasonable to assume that the level parameters are independent of each other and each of them is independent of the error terms  $\varepsilon_{ijk}$ . To avoid confusion, we summarise the dimension variables as follows:

$$\begin{aligned}
 p &= \text{number of fixed effects,} \\
 q^{\text{RL3}} &= \text{number of random effects at the group level,} \\
 q^{\text{RL2}} &= \text{number of random effects at the subgroup level,} \\
 q^{\text{G}} &= \text{number of spline basis functions,} \\
 m &= \text{total number of groups} \\
 n_i &= \text{number of subgroups in the } i\text{th group,} \\
 o_{ij} &= \text{number of observations in the } j\text{th subgroup within the } i\text{th group,} \\
 N &= \text{total number of subgroups, i.e. } \sum_{i=1}^m n_i, \\
 O &= \text{total number of observations, i.e. } \sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}.
 \end{aligned}$$

Much like in the two-level case, model (3.7) can be easily extended to other response distributions and including semiparametric extensions and hence we omit the details here.

#### 3.7.1 Mixed model representation

In keeping with the notation used in Section 3.3, we express model (3.7) in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} = \mathbf{X}^{\text{R}}\boldsymbol{\beta}^{\text{R}} + \mathbf{X}^{\text{G}}\boldsymbol{\beta}^{\text{G}} + \mathbf{Z}^{\text{R}}\mathbf{u}^{\text{R}} + \mathbf{Z}^{\text{G}}\mathbf{u}^{\text{G}} \quad (3.16)$$

by partitioning the following vectors and matrices:

$$\mathbf{X}^{\text{R}} \equiv \begin{bmatrix} \mathbf{X}_1^{\text{R}} \\ \vdots \\ \mathbf{X}_m^{\text{R}} \end{bmatrix}, \quad \mathbf{X}^{\text{G}} \equiv \begin{bmatrix} \mathbf{X}_1^{\text{G}} \\ \vdots \\ \mathbf{X}_m^{\text{G}} \end{bmatrix}, \quad \mathbf{Z}^{\text{R}} \equiv \text{blockdiag}_{1 \leq i \leq m} \left( \mathbf{X}_i^{\text{R}} \mid \text{blockdiag}_{1 \leq j \leq n_i} (\mathbf{X}_{ij}^{\text{R}}) \right).$$



$$\mathbf{u}^R \equiv \begin{bmatrix} \mathbf{u}_1^R \\ \vdots \\ \mathbf{u}_m^R \end{bmatrix}, \quad \mathbf{u}^G \equiv \begin{bmatrix} \mathbf{u}_1^G \\ \vdots \\ \mathbf{u}_L^G \end{bmatrix}, \quad \mathbf{Z}^G \equiv \begin{bmatrix} z_{\ell 1}(\mathbf{X}_1^G) & \cdots & z_{\ell q_\ell^G}(\mathbf{X}_1^G) \\ \vdots & \ddots & \vdots \\ z_{\ell 1}(\mathbf{X}_m^G) & \cdots & z_{\ell q_\ell^G}(\mathbf{X}_m^G) \end{bmatrix},$$

where

$$\mathbf{X}_i^R \equiv \begin{bmatrix} \mathbf{X}_{i1}^R \\ \vdots \\ \mathbf{X}_{in_m}^R \end{bmatrix} \quad \text{given} \quad \mathbf{X}_{ij}^R \equiv \begin{bmatrix} 1 & x_{ij1} \\ \vdots & \vdots \\ 1 & x_{ij o_{mnm}} \end{bmatrix}, \quad \mathbf{X}_i^G \equiv \begin{bmatrix} \mathbf{X}_{i1}^C \\ \vdots \\ \mathbf{X}_{in_m}^C \end{bmatrix}$$

$$\mathbf{u}_i^R \equiv \begin{bmatrix} \mathbf{u}_i^{\text{RL3}} \\ \mathbf{u}_{i1}^{\text{RL2}} \\ \vdots \\ \mathbf{u}_{in_i}^{\text{RL2}} \end{bmatrix} \quad \text{given} \quad \mathbf{u}_i^{\text{RL3}} \equiv \begin{bmatrix} u_{0i}^{\text{RL3}} \\ u_{1i}^{\text{RL3}} \end{bmatrix} \quad \text{and} \quad \mathbf{u}_{ij}^{\text{RL2}} \equiv \begin{bmatrix} u_{0ij}^{\text{RL2}} \\ u_{1ij}^{\text{RL2}} \end{bmatrix}, \quad \mathbf{u}_\ell^G \equiv \begin{bmatrix} u_\ell^G \\ \vdots \\ u_{q_\ell^G}^G \end{bmatrix}.$$

The response variable  $\mathbf{y}$  is an  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times 1$  vector and  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of fixed effects. The matrices  $\mathbf{X}^R$  and  $\mathbf{Z}^R$  are the  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times q^{\text{RL3}}$  and  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times \{m(q^{\text{RL2}} + q^{\text{RL3}} \sum_{i=1}^m n_i)\}$  general design matrices corresponding to the fixed effects vector  $\boldsymbol{\beta}^R$  and random group effects  $\mathbf{u}^R$  respectively. The matrices  $\mathbf{X}^G$  and  $\mathbf{Z}^G$  are the  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times 1$  and  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times \sum_{\ell=1}^L q_\ell^G$  general design matrices corresponding to the fixed effects vector  $\boldsymbol{\beta}^G$  and spline coefficients  $\mathbf{u}^G$  respectively. The matrices  $\boldsymbol{\Sigma}^{\text{RL2}}$  and  $\boldsymbol{\Sigma}^{\text{RL3}}$  are the  $q^{\text{RL2}} \times q^{\text{RL2}}$  and  $q^{\text{RL3}} \times q^{\text{RL3}}$  *unstructured* (no pattern is specified) covariance matrices. The random effects covariance matrix of  $(\mathbf{u}^R, \mathbf{u}^G)$  is

$$\text{Cov} \left( \begin{bmatrix} \mathbf{u}^R \\ \mathbf{u}^G \end{bmatrix} \right) = \begin{bmatrix} \text{blockdiag}_{1 \leq i \leq m} \begin{bmatrix} \boldsymbol{\Sigma}^{\text{RL3}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \boldsymbol{\Sigma}^{\text{RL2}} \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} (\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) \end{bmatrix}.$$

Here the  $\mathbf{u}_\ell^G$ s are  $q_\ell^G \times 1$  vectors for  $1 \leq \ell \leq L$ , so the blocks of  $\text{Cov}(\mathbf{u}^G)$  corresponding to the decomposition of  $\mathbf{u}^G$ .

We are now ready to present the full three-level Bayesian semiparametric mixed models. Figure 3.5 is the directed acyclic graph corresponding to model (3.17). We use the same prior specification for the model parameters as described in Section 2.4. Consider the response vector  $\mathbf{y}$  with either the Gaussian, Student- $t$ , Bernoulli or Poisson distribution as follows:

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 \sim \text{Dist}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \mathbf{R}), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), \quad (3.17)$$

### 3.7. EXTENSION TO THREE-LEVEL SEMIPARAMETRIC MIXED MODELS

$$\begin{aligned}
 \mathbf{u} | \sigma_{ul}^2, \Sigma^{\text{RL2}}, \Sigma^{\text{RL3}} &\sim N \left( \mathbf{0}, \begin{bmatrix} \text{blockdiag} \left( \sigma_{ul}^2 \mathbf{I}_{q_\ell^G} \right) & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left[ \begin{matrix} \Sigma^{\text{RL3}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \Sigma^{\text{RL2}} \end{matrix} \right] \end{bmatrix} \right), \\
 \Sigma^{\text{RL2}} | a_1^{\text{RL2}}, \dots, a_{q^{\text{RL2}}}^{\text{RL2}} &\sim \text{Inverse-Wishart} \left( \nu + q^{\text{RL2}} - 1, 2\nu \text{diag}(1/a_1^{\text{RL2}}, \dots, 1/a_{q^{\text{RL2}}}^{\text{RL2}}) \right), \\
 \Sigma^{\text{RL3}} | a_1^{\text{RL3}}, \dots, a_{q^{\text{RL3}}}^{\text{RL3}} &\sim \text{Inverse-Wishart} \left( \nu + q^{\text{RL3}} - 1, 2\nu \text{diag}(1/a_1^{\text{RL3}}, \dots, 1/a_{q^{\text{RL3}}}^{\text{RL3}}) \right), \\
 a_{r_{L2}}^{\text{RL2}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{\text{RL2}}^2 \right), \quad r_{L2} = 1, \dots, q^{\text{RL2}}, \\
 a_{r_{L3}}^{\text{RL3}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{\text{RL3}}^2 \right), \quad r_{L3} = 1, \dots, q^{\text{RL3}}, \\
 \sigma_{ul}^2 | a_{ul} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/a_{ul} \right), \quad a_{ul} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{ul}^2 \right),
 \end{aligned}$$

where  $a_{r_{L2}}^{\text{RL2}}$  ( $1 \leq r_{L2} \leq q^{\text{RL2}}$ ),  $a_{r_{L3}}^{\text{RL3}}$  ( $1 \leq r_{L3} \leq q^{\text{RL3}}$ ) and  $a_{ul}$  ( $1 \leq \ell \leq L$ ) are the auxiliary variables used for the scale parameters in the model. The notation ‘‘Dist’’ denotes the type of response distribution and  $\mathbf{R}$  represents the additional parameters corresponding to the errors, e.g.  $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}$  for the Gaussian response mixed model. Fitting (3.17) involves setting of the hyperparameters  $\sigma_\beta^2$ ,  $A_{\text{RL2}}$ ,  $A_{\text{RL3}}$  and  $A_{ul}$  by the analyst. To ensure the property of Huang and Wand (2013) holds, we set  $\nu = 2$ .

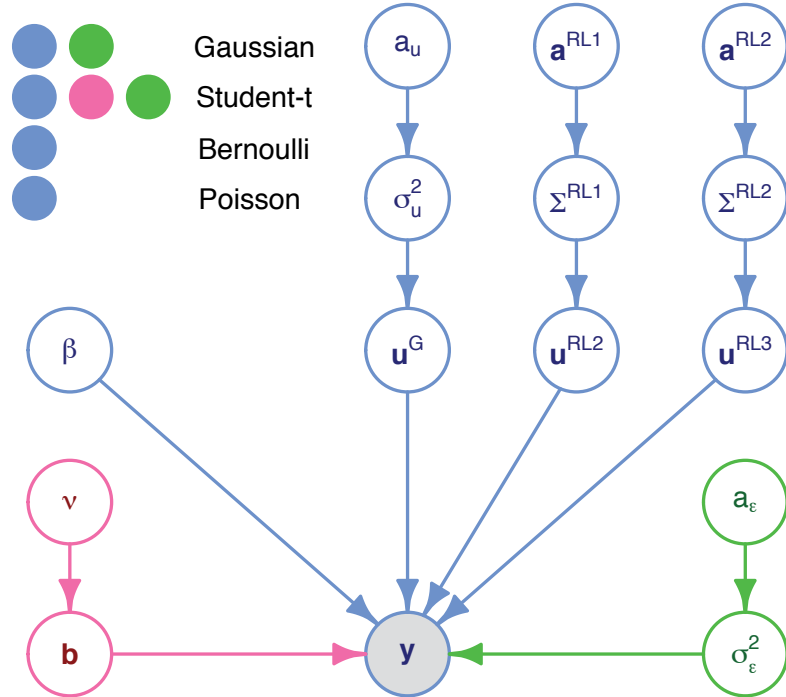


Figure 3.5: Directed acyclic graphs for model (3.17). The shaded node corresponds to the observed data vector. The colour keys at the top of the figure indicate the components of the graph corresponding to each response outcome.

### 3.7.2 Naïve mean field variational Bayes algorithms

From an extendability standpoint, a general design framework for mixed models is a particularly attractive approach that allows MFVB algorithms to cater for larger and more complicated models with straightforward algebraic adjustments. The structural changes in the random effects matrices  $\mathbf{Z}$  and  $\mathbf{G}$  underpin the complexity of a mixed model being transitioned from a two-level to a three-level hierarchical structure. For example, Figure 3.6 shows that the  $\mathbf{Z}$  matrix starts off as having a simple block-diagonal structure and, as the level of hierarchy increases, it grows into a “nested” block-diagonal structure (each large block is itself block-diagonal, with one or more small blocks on the diagonal). Further, even with the additional random group effects parameters, the covariance matrix shown in (3.17) remains to be a block-diagonal, but certainly larger, matrix with different entries for the variance components corresponding to the random basis coefficients for the overall mean and random group effects. Taking these into consideration, it then follows that the MFVB algorithms for model (3.17) simply involves modifications of Algorithms 7, 8 and 9 to incorporate the aforementioned structural changes in the  $\mathbf{Z}$  and  $\mathbf{G}$  matrices, as well as updates for the additional model parameters as follows:

$q^*(\Sigma^{\text{RL2}})$  is the Inverse-Wishart  $\left(\nu + \sum_{i=1}^m n_i + q^{\text{RL2}} + 1, \mathbf{B}_{q(\Sigma^{\text{RL2}})}\right)$  density function,  
 $q^*(a_{\tau_{L2}}^{\text{RL2}})$  is the Inverse-Gamma  $\left(\frac{1}{2}(\nu + q^{\text{RL2}}), B_{q(a_{\tau_{L2}}^{\text{RL2}})}\right)$  density function,  
 $q^*(\Sigma^{\text{RL3}})$  is the Inverse-Wishart  $\left(\nu + m + q^{\text{RL3}} + 1, \mathbf{B}_{q(\Sigma^{\text{RL3}})}\right)$  density function and  
 $q^*(a_{\tau_{L3}}^{\text{RL3}})$  is the Inverse-Gamma  $\left(\frac{1}{2}(\nu + q^{\text{RL3}}), B_{q(a_{\tau_{L3}}^{\text{RL3}})}\right)$  density function.

These naïve mean field variational Bayes algorithms are summarised in Algorithm 10.

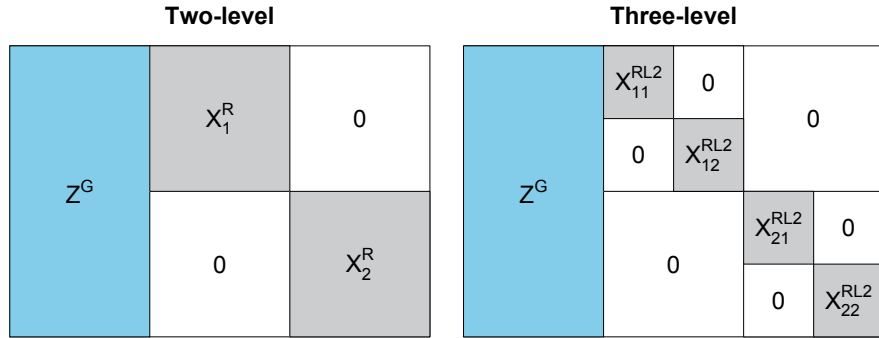


Figure 3.6: Different  $\mathbf{Z}$  structures for the two-level and three-level models with  $m = 2$ . The  $\mathbf{Z}$  matrix starts off as having a simple block-diagonal structure and, as the level of hierarchy increases, it grows into a “nested” block-diagonal structure. The definitions of matrices are described in the text.

### 3.7.3 Streamlined mean field variational Bayes algorithms

As pointed out in Section 3.4, the naïve implementation of MFVB algorithms for arbitrarily large grouped data is extremely inefficient in terms of speed and storage. The time complexity can be as high as  $O(m^3)$ , making variational inference impractical for large and complex mixed models. Through exploiting the inherent block-diagonal structure of the effects covariance matrix, we develop fast and memory-efficient MFVB algorithms that streamline inversion and update for  $\Sigma_{q(\beta, \mathbf{u})}$ . Here we extend the streamlined approach to the three-level model settings. The essential trick for this extension is to permute the  $\Sigma_{q(\beta, \mathbf{u})}$  matrix into an approximate block-diagonal form for decomposition, as shown in Section 3.4. Then every previously described step can be analogously applied to the three-level mixed model by simply modifying the submatrices  $\mathbf{F}$ ,  $\mathbf{G}_i$  and  $\mathbf{H}_i$  that are defined in Section 3.4. Take the Gaussian response model as an example,

$$\begin{aligned} \mathbf{F} &\equiv \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}^G)^\top \mathbf{C}^G + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} \left( \mu_{q(1/\sigma_{u_\ell}^2)} \mathbf{I}_{q_\ell^G} \right) \end{bmatrix}, \\ \mathbf{G}_i &\equiv \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}_i^G)^\top \mathbf{Z}_i^R, \\ \text{and } \mathbf{H}_i &\equiv \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{Z}_i^R)^\top \mathbf{Z}_i^R + \begin{bmatrix} \Sigma^{\text{RL3}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \Sigma^{\text{RL2}} \end{bmatrix} \right\}^{-1}. \end{aligned}$$

Despite the significant amount of time we have already saved, unlike the two-level case, the resultant block diagonal structure in  $\mathbf{H}_i$  can still be considered large when  $n_i$  is large and comprises unnecessary zeros. This hinders the optimal time savings that would be possible and further research into this area is undergoing.

## 3.8 Concluding remarks

Mean field variational Bayes with streamlining, as exemplified by Algorithms 7, 8 and 9, is a useful addition to the longitudinal and multilevel data analysis arsenal. It allows rapid and accurate approximate Bayesian inference for very large datasets with computational times that increase only linearly in the number of groups. Extension of streamlined MFVB to higher level longitudinal and multilevel models is also of interest, although this requires a considerable amount of analytic computation. This chapter has helped opened up a new branch of methodology for fitting and inference for grouped data in the high volume/velocity era.

### 3.8. CONCLUDING REMARKS

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mathbf{M}_{q((\boldsymbol{\Sigma}_u^{\text{RL2}})^{-1})}$ , a  $q^{\text{RL2}} \times q^{\text{RL2}}$  positive definite matrix,  $\mathbf{M}_{q((\boldsymbol{\Sigma}_u^{\text{RL3}})^{-1})}$ , a  $q^{\text{RL3}} \times q^{\text{RL3}}$  positive definite matrix,  $\mu_{q(1/a_{\text{RL2}}^{\text{RL2}})} > 0$ ,  $1 \leq r_{\text{L2}} \leq q^{\text{RL2}}$ ,  $\mu_{q(1/a_{\text{RL3}}^{\text{RL3}})} > 0$ ,  $1 \leq r_{\text{L3}} \leq q^{\text{RL3}}$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $1 \leq \ell \leq L$ ,  $\boldsymbol{\xi}$ , a  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times 1$  vector of positive entries and  $\mu_{q(1/b_i)} > 0$ ,  $1 \leq i \leq O$ ,  $\mu_{q(\nu)} > 0$ ,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , a  $(p + m q^{\text{RL3}} + \sum_{i=1}^m n_i q^{\text{RL2}} + \sum_{\ell=1}^L q_\ell^G) \times 1$  vector and  $\mathbf{w}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , a  $(\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}) \times 1$  vector of positive entries.

**Cycle through updates:**

$$\text{Define: } M_{q(\Omega)} = \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq L} (\mu_{q(1/\sigma_{u\ell}^2)} \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag}_{1 \leq i \leq m} \begin{bmatrix} \mathbf{M}_{q((\boldsymbol{\Sigma}_u^{\text{RL2}})^{-1})} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}_u^{\text{RL3}})^{-1})} \end{bmatrix} \end{bmatrix}$$

**If fitting the Gaussian model:**

- Update multivariate normal  $q^*(\boldsymbol{\beta}, \mathbf{u})$  parameters:

$$\begin{aligned} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow (\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} + M_{q(\Omega)})^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{y} \end{aligned}$$

- Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  and inverse-Gamma  $q^*(a_\varepsilon)$  parameters:

$$\begin{aligned} B_{q(\sigma_\varepsilon^2)} &\leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left\{ \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\} \\ \mu_{q(1/\sigma_\varepsilon^2)} &\leftarrow \frac{1}{2} (\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij} + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}) \end{aligned}$$

**If fitting the Student- $t$  model:**

$$O \leftarrow \sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij} \quad ; \quad \mathbf{C}_{12} \leftarrow \text{diag}_{1 \leq i \leq O} (\mu_{q(1/b_i)})$$

- Update multivariate normal  $q^*(\boldsymbol{\beta}, \mathbf{u})$  parameters:

$$\begin{aligned} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow (\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} + M_{q(\Omega)})^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{C}_{12} \mathbf{y} \end{aligned}$$

- Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  and inverse-Gamma  $q^*(a_\varepsilon)$  parameters:

$$\begin{aligned} B_{q(\sigma_\varepsilon^2)} &\leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top \mathbf{C}_{12} (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) + \text{tr}(\mathbf{C}^\top \mathbf{C}_{12} \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\} \\ \mu_{q(1/\sigma_\varepsilon^2)} &\leftarrow \frac{1}{2} (\sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij} + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}) \end{aligned}$$

- Update inverse-Gamma  $q^*(b_i)$  parameters:

For  $1 \leq i \leq O$ :

$$\begin{aligned} B_{q(b_i)} &\leftarrow \frac{1}{2} \left[ \mu_{q(\nu)} + \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})_i^2 + (\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top)_{ii} \right\} \right] \\ \mu_{q(1/b_i)} &\leftarrow \frac{1}{2} (\mu_{q(\nu)} + 1) / B_{q(b_i)} \\ \mu_{q(\log b_i)} &\leftarrow \log(B_{q(b_i)}) - \text{digamma} \left\{ \frac{1}{2} (\mu_{q(\nu)} + 1) \right\} \end{aligned}$$

- Update  $q^*(\nu)$  parameters:

$$C_1 \leftarrow \sum_{i=1}^O (\mu_{q(\log(b_i))} + \mu_{q(1/b_i)})$$

$$\mu_{q(\nu)} \leftarrow \exp \left\{ \log \mathcal{F}(1, \sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}, C_1, \nu_{\min}, \nu_{\max}) \right. \\ \left. - \log \mathcal{F}(0, \sum_{i=1}^m \sum_{j=1}^{n_i} o_{ij}, C_1, \nu_{\min}, \nu_{\max}) \right\}$$

**If fitting the Bernoulli model:**

- Update Multivariate Normal  $q^*(\beta, \mathbf{u})$  and  $\xi$  parameters:

$$\Sigma_{q(\beta, \mathbf{u})} \leftarrow [2\mathbf{C}^\top \text{diag}\{A(\xi)\}\mathbf{C} + \mathbf{M}_{q(\Omega)}]^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u})} \leftarrow \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top (\mathbf{y} - \frac{1}{2}\mathbf{1})$$

$$\xi \leftarrow \sqrt{\text{diagonal} \left\{ \mathbf{C} \left( \Sigma_{q(\beta, \mathbf{u})} + \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}^\top \right) \mathbf{C}^\top \right\}}$$

**If fitting the Poisson model:**

- Update  $q^*(\beta, \mathbf{u})$  and  $\mathbf{w}_{q(\beta, \mathbf{u})}$  parameters:

$$\Sigma_{q(\beta, \mathbf{u})} \leftarrow \left\{ \mathbf{C}^\top \text{diag}(\mathbf{w}_{q(\beta, \mathbf{u})}) \mathbf{C} + \mathbf{M}_{q(\Omega)} \right\}^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u})} \leftarrow \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \Sigma_{q(\beta, \mathbf{u})} \left\{ \mathbf{C}^\top (\mathbf{y} - \mathbf{w}_{q(\beta, \mathbf{u})}) - \mathbf{M}_{q(\Omega)} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \right\}$$

$$\mathbf{w}_{q(\beta, \mathbf{u})} \leftarrow \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top) \right\}$$

**For any response model:**

- Update inverse-Gamma  $q^*(a_{u\ell})$  and inverse-Gamma  $q^*(\sigma_{u\ell}^2)$  parameters:

For  $\ell = 1, \dots, L$ :

$$\mu_{q(1/a_{u\ell})} \leftarrow 1 / (\mu_{q(1/\sigma_{u\ell}^2)} + A_{u\ell}^{-2})$$

$$\mu_{q(\sigma_{u\ell}^2)} \leftarrow \frac{q_\ell^G + 1}{\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell^G)}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}_\ell^G)}) + 2\mu_{q(1/a_{u\ell})}}$$

- Update inverse-Gamma  $q^*(a_{r_{L2}}^{\text{RL2}})$  and inverse-Wishart  $q^*(\Sigma^{\text{RL2}})$  parameters:

For  $r_{L2} = 1, \dots, q^{\text{RL2}}$ :

$$B_{q(a_{r_{L2}}^{\text{RL2}})} \leftarrow \nu (\mathbf{M}_{q((\Sigma_u^{\text{RL2}})^{-1})})_{r_{L2} r_{L2}} + A_{\text{RL2}}^{-2}$$

$$\mu_{q(a_{r_{L2}}^{\text{RL2}})} \leftarrow \frac{1}{2} (\nu + q^{\text{RL2}}) / B_{q(a_{r_{L2}}^{\text{RL2}})}$$

$$\mathbf{B}_{q(\Sigma^{\text{RL2}})} \leftarrow \sum_{i=1}^m \sum_{j=1}^{n_i} (\boldsymbol{\mu}_{q(\mathbf{u}_{ij}^{\text{RL2}})} \boldsymbol{\mu}_{q(\mathbf{u}_{ij}^{\text{RL2}})}^\top + \Sigma_{q(\mathbf{u}_{ij}^{\text{RL2}})}) \\ + 2\nu \text{diag}(\mu_{q(1/a_1^{\text{RL2}})}, \dots, \mu_{q(1/a_{q^{\text{RL2}}}^{\text{RL2}})})$$

$$\mathbf{M}_{q((\Sigma_u^{\text{RL2}})^{-1})} \leftarrow (\nu + \sum_{i=1}^m n_i + q^{\text{RL2}} - 1) \mathbf{B}_{q(\Sigma^{\text{RL2}})}^{-1}$$

- Update inverse-Gamma  $q^*(a_{r_{L3}}^{\text{RL3}})$  and inverse-Wishart  $q^*(\Sigma^{\text{RL3}})$  parameters:

For  $r_{L3} = 1, \dots, q^{\text{RL3}}$ :

$$B_{q(a_{r_{L3}}^{\text{RL3}})} \leftarrow \nu(\mathbf{M}_{q((\Sigma_u^{\text{RL3}})^{-1})})_{r_{L3}r_{L3}} + A_{\text{RL3}}^{-2}$$

$$\mu_{q(a_{r_{L3}}^{\text{RL3}})} \leftarrow \frac{1}{2}(\nu + q^{\text{RL3}})/B_{q(a_{r_{L3}}^{\text{RL3}})}$$

$$\begin{aligned} \mathbf{B}_{q(\Sigma^{\text{RL3}})} \leftarrow & \sum_{i=1}^m (\boldsymbol{\mu}_{q(u_i^{\text{RL3}})} \boldsymbol{\mu}_{q(u_i^{\text{RL3}})}^\top + \Sigma_{q(u_i^{\text{RL3}})}) \\ & + 2\nu \text{diag}(\mu_{q(1/a_1^{\text{RL3}})}, \dots, \mu_{q(1/a_q^{\text{RL3}})}) \end{aligned}$$

$$\mathbf{M}_{q((\Sigma_u^{\text{RL3}})^{-1})} \leftarrow (\nu + m + q^{\text{RL3}} - 1) \mathbf{B}_{q(\Sigma^{\text{RL3}})}^{-1}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**

---

Algorithm 10: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the three-level Bayesian semiparametric mixed models with Gaussian, Student- $t$ , Bernoulli and Poisson responses.

## Chapter 4

# Measurement Error, Missing Data and Real-Time Extensions

*You can use all the quantitative data you can get, but you still have to distrust it and use your own intelligence and judgment.*

Alvin Toffler

### 4.1 Introduction

The beauty of the Bayesian paradigm combined with MFVB algorithms is its tremendous flexibility. Our MFVB algorithms are “modular” in the sense that the methods of handling, for example, predictor effects, spline or wavelet regression, measurement error and/or missing data models and real-time models can be combined relatively easily by incorporating additional nodes in the DAGs. In this chapter, we explore three interesting and challenging extensions of Bayesian semiparametric mixed models with an emphasis on the Gaussian response as previously described in Section 2.5. These include: data subject to measurement error problems, data subject to missing data problems and data that are processed in real-time as they arrive.

All of the examples given so far in this thesis have been for datasets with complete information on all variables and with the assumption that the variables of interest have been measured accurately. Unfortunately, many real-world applications suffer from measurement error and/or missing data problems. For example, in the *New South Wales 45 and*

---

The real-time content of this chapter was presented at the 16th J.B. Douglas Awards organised by the Statistical Society of Australia, Sydney, Australia, 2015. (Awarded Equal First Prize for oral presentation).



*Up Study* (<http://www.45andUp.org.au>) in which the data come from questionnaires filled out by individuals aged 45 years and over, it is common for questions involving smoking and alcohol consumption to be skipped or under-reported. In addition, a potentially important variable, such as weekly protein intake, is not measured by the questionnaires but a surrogate, such as the number of times eating red meats per week, could be formed from a composite of the questionnaire responses. Such a situation is often called the *classical measurement error problem*, in which the truth is measured with additive errors, usually with a constant variance. Methods on how to best account for such impurities in the data have spawned major areas of statistical research. Books devoted to handling measurement error and/or missing data problems include Little and Rubin (2014), Carroll *et al.* (2006) and Daniels and Hogan (2008).

When the data are susceptible to measurement error and/or missingness problems, a Bayesian approach allows a relatively straightforward incorporation of standard measurement error and/or missing data mechanisms (e.g. Faes *et al.*, 2011; Richardson *et al.*, 2002), resulting in a larger Bayesian hierarchical model. Inference via MCMC is simple in principle, but can be costly in computational time. The first two-third of this chapter is concerned with the fast variational analysis of Bayesian semiparametric mixed models with error-contaminated data. Efficiency is achieved by using MFVB posterior approximations. This is a deterministic approach that yields approximate inference, rather than asymptotically exact inference produced by the standard MCMC methods.

The remaining one-third of this chapter is to extend our previously developed streamlined MFVB algorithms to handle real-time variational inference, in which the data are processed as they are collected and made immediately available via modern technologies. Previous variational inference assumes that the data are processed in *batch*, that is, all at the same time. Here we treat data as arriving online and process them in real time, and thus have the advantage of not storing potentially very large datasets.

## 4.2 Gaussian semiparametric mixed models with measurement error problems

We distinguish between two types of predictors:  $\mathbf{s}$  represents a predictor that, for all practical purposes, is measured without error (for the present example,  $\mathbf{s}$  is continuous with non-linear predictor effects), and  $\mathbf{x}$  represents a predictor that is not measured accurately on all subjects. The distinguishing feature of a measurement error model from model (2.11) is that we can observe an imperfect surrogate of  $\mathbf{x}$ , denoted by  $\mathbf{w}$ , and the classical measurement error model states that  $\mathbf{w} = \mathbf{x} + \mathbf{v}$ , where  $\mathbf{w}$  is an unbiased measure of  $\mathbf{x}$  and  $\mathbf{v}$  are i.i.d.  $N(\mathbf{0}, \sigma_w^2 \mathbf{I})$  random variables. Since  $\mathbf{x}$  is not fully observed, it is not

#### 4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

possible to regress  $y$  on  $(\mathbf{x}, \mathbf{s})$ . The goal of measurement error modelling is therefore to obtain accurate estimates of model parameters indirectly by fitting a model of  $\mathbf{y}$  in terms of  $(\mathbf{w}, \mathbf{s})$ .

Early work on Bayesian analysis of measurement error problems has been developed by Clayton (1992), followed by Richardson *et al.* (2002). Their work are based on the notion of *structural and functional specifications*. We split the predictor  $\mathbf{x}$  into two subgroups:  $\mathbf{x}_{\text{obs}}$  that involves the accurate measurements of the predictor and  $\mathbf{x}_{\text{unobs}}$  that involves the inaccurate measurements of the predictor. The observations  $\mathbf{x}_{\text{obs}}$  are generally obtained from a *validation sample*, i.e. a subsample of the data where, by design, the predictor  $\mathbf{x}$  is accurately recorded by the use of the so-called “gold standard” method. Quite often, this is costly to implement on a large scale, hence  $\mathbf{x}_{\text{unobs}}$  is usually larger than  $\mathbf{x}_{\text{obs}}$ . To be specific, let us distinguish data in the validation sample and the main sample by the superscript index “obs” or “unobs” respectively. Hence,  $\mathbf{w}_{\mathbf{x}_{\text{obs}}}$  partners the validation sample  $\mathbf{x}_{\text{obs}}$  and  $\mathbf{w}_{\mathbf{x}_{\text{unobs}}}$  contains the surrogates for the  $\mathbf{x}_{\text{obs}}$ .

Throughout this section, we let  $\mathbf{y}$  denote the known outcome,  $\mathbf{x}$  denote the true predictor which is unobserved (except in the validation sample),  $\mathbf{w}$  denote the observed surrogate for  $\mathbf{x}$  and  $\mathbf{s}$  denote the known continuous, non-linear predictor. For the present work, we assume that  $\mathbf{x}$  and  $\mathbf{w}$  are univariate. As detailed in Richardson *et al.* (2002), the structural specifications of a measurement error model entail the formulation of three submodels:

1. a regression model which expresses the relationship between  $\mathbf{y}$ ,  $\mathbf{x}$  and  $\mathbf{s}$ , denoted by  $p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \mathbf{s}, \mathbf{u}, \sigma_\varepsilon^2)$ ;
2. a classical measurement model which express the relationship between  $\mathbf{w}$  and  $\mathbf{x}$ , denoted by  $p(\mathbf{w}|\mathbf{x}, \sigma_w^2)$ ; and
3. a prior model which specifies the prior distribution of  $\mathbf{x}$ , denoted by  $p(\mathbf{x}|\mu_x, \sigma_x^2)$ ,

The definition of  $\mathbf{u}$  is given in Section 2.3 and appropriate priors are placed on  $\sigma_\varepsilon^2, \sigma_w^2, \mu_x$  and  $\sigma_x^2$ . Following the above, functional forms for the distributions involved in the submodels are required to be chosen. For the present regression and measurement error models, we use linear regression coupled with normal error models corresponding to additive errors, a choice that has been frequently used in epidemiology. To keep this section manageable, we also choose a normal distribution to represent the prior model  $p(\mathbf{x}|\mu_x, \sigma_x^2)$ . Although not the most ideal, this can be relatively straightforward to extend to a mixture model, which avoids possible model misspecification in the prior model.

Putting all of these components together, we now consider the following model which exploits the structure described above:

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \mathbf{X}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2\mathbf{I}), & \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_\beta^2\mathbf{I}) \\ \mathbf{x}_{\text{obs}}|\mu_x, \sigma_x^2 &\sim N(\mu_x\mathbf{1}_{n_{\text{obs}}}, \sigma_x^2\mathbf{I}), & \mathbf{x}_{\text{unobs}}|\mu_x, \sigma_x^2 &\sim N(\mu_x\mathbf{1}_{n_{\text{unobs}}}, \sigma_x^2\mathbf{I}), \end{aligned} \quad (4.1)$$

#### 4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

$$\begin{aligned}
\mathbf{w}_{x_{\text{obs}}} | \mathbf{x}_{\text{obs}}, \sigma_w^2 &\sim N(\mathbf{x}_{\text{obs}}, \sigma_w^2 \mathbf{I}), & \mathbf{w}_{x_{\text{unobs}}} | \mathbf{x}_{\text{unobs}}, \sigma_w^2 &\sim N(\mathbf{x}_{\text{unobs}}, \sigma_w^2 \mathbf{I}) \\
\mathbf{u} | \sigma_u^2, \Sigma^{\text{R}} &\sim N\left(\mathbf{0}, \begin{bmatrix} \sigma_u^2 \mathbf{I}_{q^{\text{G}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \Sigma^{\text{R}} \end{bmatrix}\right), & \mu_x | \sigma_x^2 &\sim N(0, \sigma_x^2), \\
\sigma_x^2 | a_x &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_x), & a_x &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_x^2), \\
\Sigma^{\text{R}} | a_1^{\text{R}}, \dots, a_{q^{\text{R}}}^{\text{R}} &\sim \text{Inverse-Wishart}\left(\nu + q^{\text{R}} - 1, 2\nu \text{diag}(1/a_1^{\text{R}}, \dots, 1/a_{q^{\text{R}}}^{\text{R}})\right), \\
a_r^{\text{R}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_r^2), & 1 \leq r \leq q^{\text{R}}, \\
\sigma_u^2 | a_u &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_u), & a_u &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_u^2), \\
\sigma_w^2 | a_w &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_w), & a_w &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_w^2), \\
\sigma_\varepsilon^2 | a_\varepsilon &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_\varepsilon), & a_\varepsilon &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2).
\end{aligned}$$

Whilst model (4.1) extends model (2.11) by including the measurement error mechanism, it follows the same way to be represented as a mixed model as shown in Section 2.3. The prior specification of model (4.1) can be found in Section 2.4. The conditional independence properties of model (4.1) can be visualised in the DAG as depicted in Figure 4.2.

##### 4.2.1 Notation

The following notation is useful in the upcoming subsections. Let  $n_i^{\text{obs}}$  denote the number of observed  $x_{ij}$ s in the  $i$ th group and  $n_i^{\text{unobs}}$  denote the number of unobserved  $x_{ij}$ s in the  $i$ th group. Let  $\mathbf{x}_{\text{obs},i}$  be the  $n_i^{\text{obs}} \times 1$  vector containing the observed  $x_{ij}$ s and  $\mathbf{x}_{\text{unobs},i}$  be the  $n_i^{\text{unobs}} \times 1$  containing the unobserved  $x_{ij}$ s. We reorder the data so that the observed data is first within each group. Hence, the full vector of predictor is

$$\mathbf{x} \equiv [\mathbf{x}_{\text{obs},1} \quad \mathbf{x}_{\text{unobs},1} \quad \cdots \quad \mathbf{x}_{\text{obs},m} \quad \mathbf{x}_{\text{unobs},m}]^{\text{T}}.$$

In addition, let  $x_{\text{obs},ij}$  be the value of the predictor corresponding to the  $j$ th entry of  $\mathbf{x}_{\text{obs},i}$  and  $x_{\text{unobs},ij'}$  be the value of the predictor corresponding to the  $j'$ th entry of  $\mathbf{x}_{\text{unobs},i}$  for  $1 \leq i \leq m$ ,  $1 \leq j \leq n_i^{\text{obs}}$  and  $1 \leq j' \leq n_i^{\text{unobs}}$ .

Define the following vectors and matrices

$$\begin{aligned}
\mathbf{x}_i &= \begin{bmatrix} \mathbf{x}_{\text{obs},i} \\ \mathbf{x}_{\text{unobs},i} \end{bmatrix}, & \mathbf{x}_{\text{obs},i} &= \begin{bmatrix} x_{\text{obs},i1} \\ \vdots \\ x_{\text{obs},in_i^{\text{obs}}} \end{bmatrix}, & \mathbf{x}_{\text{unobs},i} &= \begin{bmatrix} x_{\text{unobs},i1} \\ \vdots \\ x_{\text{unobs},in_i^{\text{unobs}}} \end{bmatrix}, \\
\mathbf{Z}_{x_{\text{obs}}}^{\text{R}} &= \text{blockdiag}\left(\left[\mathbf{1} \quad \mathbf{x}_{\text{obs},i}\right]\right)_{1 \leq i \leq m}, & \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} &= \text{blockdiag}\left(\left[\mathbf{1} \quad \mathbf{x}_{\text{unobs},i}\right]\right)_{1 \leq i \leq m},
\end{aligned}$$

4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

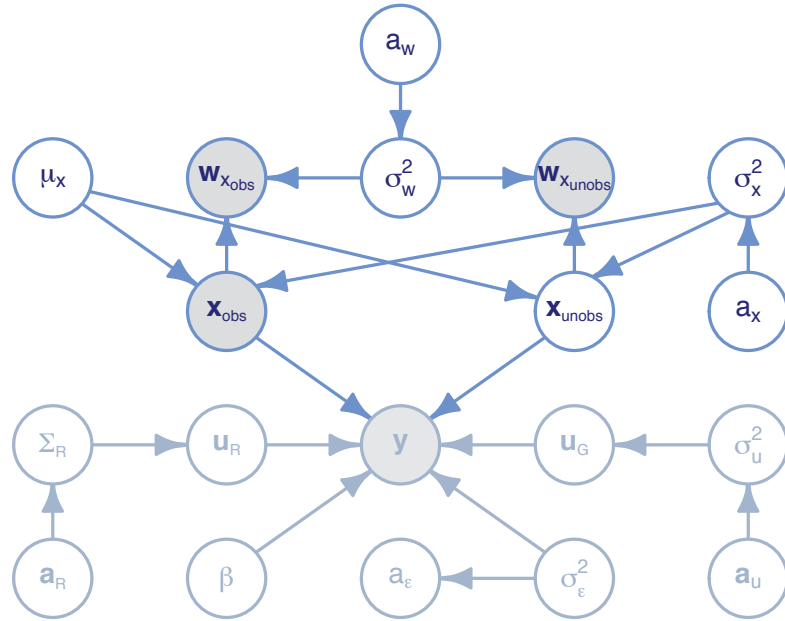


Figure 4.1: Directed acyclic graph corresponds to model (4.1). The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are presented in the Gaussian model DAG as shown in Figure 2.3.

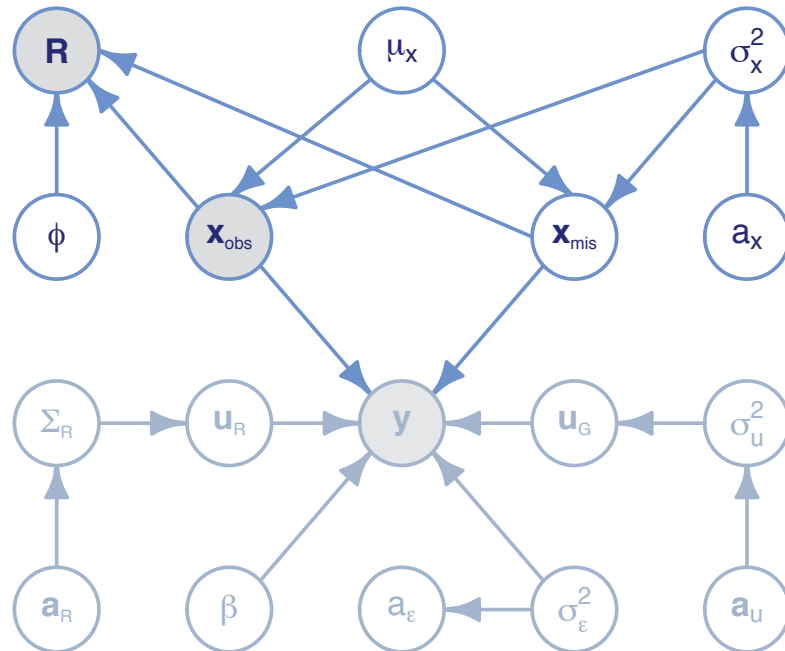


Figure 4.2: Directed acyclic graph corresponds to model (4.7). The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are presented in the Gaussian model DAG as shown in Figure 2.3.

4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

$$\mathbf{Z}_{x_{\text{obs}}}^{\text{G}} = \begin{bmatrix} z_1(\mathbf{s}_{x_{\text{obs},1}}) & \cdots & z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{obs},1}}) \\ \vdots & & \vdots \\ z_1(\mathbf{s}_{x_{\text{obs},m}}) & \cdots & z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{obs},m}}) \end{bmatrix}, \quad \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} = \begin{bmatrix} z_1(\mathbf{s}_{x_{\text{unobs},1}}) & \cdots & z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{unobs},1}}) \\ \vdots & & \vdots \\ z_1(\mathbf{s}_{x_{\text{unobs},m}}) & \cdots & z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{unobs},m}}) \end{bmatrix}.$$

For conciseness of forthcoming expressions we adopt the following notations:

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{x_{\text{obs}}} \\ \mathbf{C}_{x_{\text{unobs}}} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} \\ \mathbf{1} & \mathbf{x}_{\text{unobs}} & \mathbf{s}_{x_{\text{unobs}}} & \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} \end{bmatrix}$$

$$E_q(\mathbf{C}) = \begin{bmatrix} \mathbf{C}_{x_{\text{obs}}} \\ E_q(\mathbf{C}_{x_{\text{unobs}}}) \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} \\ \mathbf{1} & E_q(\mathbf{x}_{\text{unobs}}) & \mathbf{s}_{x_{\text{unobs}}} & \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & E_q(\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) \end{bmatrix},$$

where

$$E_q(\mathbf{x}_{\text{unobs}}) = [\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},1}) \cdots \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},m})]^\top$$

and  $E_q(\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) = \text{blockdiag}([\mathbf{1} \ \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},i})])_{1 \leq i \leq m}$ ,

and the explicit expression for  $\mathbf{C}^\top \mathbf{C}$  is

$$\mathbf{C}^\top \mathbf{C} = \begin{bmatrix} n & \mathbf{1}^\top \mathbf{x} & \mathbf{1}^\top \mathbf{s} \\ \mathbf{x}^\top \mathbf{1} & \mathbf{x}_{\text{obs}}^\top \mathbf{x}_{\text{obs}} + \mathbf{x}_{\text{unobs}}^\top \mathbf{x}_{\text{unobs}} & \mathbf{x}_{\text{obs}}^\top \mathbf{s}_{x_{\text{obs}}} + \mathbf{x}_{\text{unobs}}^\top \mathbf{s}_{x_{\text{unobs}}} \\ \mathbf{s}^\top \mathbf{1} & \mathbf{s}_{x_{\text{obs}}}^\top \mathbf{x}_{\text{obs}} + \mathbf{s}_{x_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} & \mathbf{s}_{x_{\text{obs}}}^\top \mathbf{s}_{x_{\text{obs}}} + \mathbf{s}_{x_{\text{unobs}}}^\top \mathbf{s}_{x_{\text{unobs}}} \\ (\mathbf{Z}^{\text{G}})^\top \mathbf{1} & (\mathbf{Z}_{x_{\text{obs}}}^{\text{G}})^\top \mathbf{x}_{\text{obs}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^\top \mathbf{x}_{\text{unobs}} & (\mathbf{Z}_{x_{\text{obs}}}^{\text{G}})^\top \mathbf{s}_{x_{\text{obs}}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^\top \mathbf{s}_{x_{\text{unobs}}} \\ (\mathbf{Z}^{\text{R}})^\top \mathbf{1} & (\mathbf{Z}_{x_{\text{obs}}}^{\text{R}})^\top \mathbf{x}_{\text{obs}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^\top \mathbf{x}_{\text{unobs}} & (\mathbf{Z}_{x_{\text{obs}}}^{\text{R}})^\top \mathbf{s}_{x_{\text{obs}}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^\top \mathbf{s}_{x_{\text{unobs}}} \\ \mathbf{1}^\top \mathbf{Z}^{\text{G}} & & \mathbf{1}^\top \mathbf{Z}^{\text{R}} \\ \mathbf{x}_{\text{obs}}^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} + \mathbf{x}_{\text{unobs}}^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & & \mathbf{x}_{\text{obs}}^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} + \mathbf{x}_{\text{unobs}}^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} \\ \mathbf{s}_{x_{\text{obs}}}^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} + \mathbf{s}_{x_{\text{unobs}}}^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & & \mathbf{s}_{x_{\text{obs}}}^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} + \mathbf{s}_{x_{\text{unobs}}}^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} \\ (\mathbf{Z}_{x_{\text{obs}}}^{\text{G}})^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & & (\mathbf{Z}_{x_{\text{obs}}}^{\text{G}})^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} \\ (\mathbf{Z}_{x_{\text{obs}}}^{\text{R}})^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}} & & (\mathbf{Z}_{x_{\text{obs}}}^{\text{R}})^\top \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} + (\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}} \end{bmatrix}. \quad (4.2)$$

Let  $\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},i}) = [\mathbf{1} \ \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},1})]$ , then the expressions corresponding to  $E_q(\mathbf{C}^\top \mathbf{C})$  with respect to  $\mathbf{x}_{\text{unobs}}$  are as follows:

First we note that:

$$E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^\top \mathbf{1}) = E_q^\top(\mathbf{1}^\top \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})$$

$$E_q((\mathbf{x}_{\text{unobs}})^\top \mathbf{s}_{x_{\text{unobs}}}) = E_q^\top((\mathbf{s}_{x_{\text{unobs}}})^\top \mathbf{x}_{\text{unobs}})$$

#### 4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

$$\begin{aligned} E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^{\top} \mathbf{s}_{x_{\text{unobs}}}) &= E_q^{\top}(\mathbf{s}_{x_{\text{unobs}}}^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) \\ E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^{\top} \mathbf{x}_{\text{unobs}}) &= E_q^{\top}(\mathbf{x}_{\text{unobs}}^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}}). \end{aligned}$$

The remaining expectation update expressions are derived via straightforward algebra manipulations:

$$\begin{aligned} E_q((\mathbf{x}_{\text{unobs}})^{\top} \mathbf{x}_{\text{unobs}}) &\leftarrow \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}})\|^2 + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 & (4.3) \\ E_q((\mathbf{s}_{x_{\text{unobs}}})^{\top} \mathbf{x}_{\text{unobs}}) &\leftarrow \sum_{i=1}^m \mathbf{s}_{x_{\text{unobs},i}}^{\top} \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},i}) \\ E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) &\leftarrow \text{blockdiag}_{1 \leq i \leq m} \left( \boldsymbol{\mu}_q^{\top}(\mathbf{X}_{\text{unobs},i}) \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},i}) \right) + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \\ E_q(\mathbf{1}^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) &\leftarrow \left[ \mathbf{1}^{\top} \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},1}) \quad \cdots \quad \mathbf{1}^{\top} \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},m}) \right] \\ E_q(\mathbf{x}_{\text{unobs}}^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) &\leftarrow \left[ \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},1})\|^2 \quad \cdots \quad \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},m})\|^2 \right] + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \\ E_q(\mathbf{s}_{x_{\text{unobs}}}^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) &\leftarrow \left[ \mathbf{s}_{x_{\text{unobs},1}}^{\top} \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},1}) \quad \cdots \quad \mathbf{s}_{x_{\text{unobs},m}}^{\top} \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},m}) \right] \\ E_q((\mathbf{x}_{\text{unobs}})^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}}) &\leftarrow \left[ \sum_{i=1}^m \boldsymbol{\mu}_q^{\top}(\mathbf{x}_{\text{unobs},i}) z_{1}(\mathbf{s}_{x_{\text{unobs},i}}) \quad \cdots \quad \sum_{i=1}^m \boldsymbol{\mu}_q^{\top}(\mathbf{x}_{\text{unobs},i}) z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{unobs},i}}) \right] \\ E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{R}})^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{G}}) &\leftarrow \begin{bmatrix} \boldsymbol{\mu}_q^{\top}(\mathbf{X}_{\text{unobs},1}) z_{1}(\mathbf{s}_{x_{\text{unobs},1}}) & \cdots & \boldsymbol{\mu}_q^{\top}(\mathbf{X}_{\text{unobs},1}) z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{unobs},1}}) \\ \vdots & & \vdots \\ \boldsymbol{\mu}_q^{\top}(\mathbf{X}_{\text{unobs},m}) z_{1}(\mathbf{s}_{x_{\text{unobs},m}}) & \cdots & \boldsymbol{\mu}_q^{\top}(\mathbf{X}_{\text{unobs},m}) z_{q^{\text{G}}}(\mathbf{s}_{x_{\text{unobs},m}}) \end{bmatrix} \\ E_q((\mathbf{Z}_{x_{\text{unobs}}}^{\text{G}})^{\top} \mathbf{Z}_{x_{\text{unobs}}}^{\text{R}}) &\leftarrow \begin{bmatrix} z_{1}^{\top}(\mathbf{s}_{x_{\text{unobs},1}}) \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},1}) & \cdots & z_{1}^{\top}(\mathbf{s}_{x_{\text{unobs},m}}) \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},1}) \\ \vdots & & \vdots \\ z_{q^{\text{G}}}^{\top}(\mathbf{s}_{x_{\text{unobs},1}}) \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},m}) & \cdots & z_{q^{\text{G}}}^{\top}(\mathbf{s}_{x_{\text{unobs},m}}) \boldsymbol{\mu}_q(\mathbf{X}_{\text{unobs},m}) \end{bmatrix}. \end{aligned}$$

#### 4.2.2 Approximate Bayesian inference via mean field variational Bayes

We now provide details on MFVB fitting and inference catered to the Gaussian response model with classical measurement errors as described in model (4.1). The essence of our MFVB approach lies in approximating the full conditional joint posterior density function of the form

$$\begin{aligned} p(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_{\varepsilon}, \mu_x, \boldsymbol{\Sigma}^{\text{R}}, \sigma_u^2, \sigma_{\varepsilon}^2, \sigma_x^2, \sigma_w^2 | \mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{w}) &\approx \\ q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_{\varepsilon}, \mu_x, \boldsymbol{\Sigma}^{\text{R}}, \sigma_u^2, \sigma_{\varepsilon}^2, \sigma_x^2, \sigma_w^2) & \end{aligned}$$

#### 4.2. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MEASUREMENT ERROR PROBLEMS

---

subject to the  $q$ -density product restriction

$$q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \mu_x, \boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_w^2) = \quad (4.4)$$

$$q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \mathbf{u}, \mathbf{a}^R, \mathbf{a}_u, a_\varepsilon, \mu_x) \times q(\boldsymbol{\Sigma}^R, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_w^2).$$

Restriction (4.4) is the minimal factorisation in the MFVB approximation, deriving based on the heuristic arguments provided by Menictas and Wand (2013). One can also use the concepts of induced factorisation and moralisation introduced in Subsection 2.5.1 to establish additional product restriction

$$q(\mathbf{x}_{\text{unobs}}) q(\boldsymbol{\beta}, \mathbf{u}) \left\{ \prod_{r=1}^{q^R} q(a_r^R) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\} q(a_\varepsilon) q(\mu_x) q(\boldsymbol{\Sigma}^R) \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \quad (4.5)$$

$$\times q(\sigma_\varepsilon^2) q(\sigma_x^2) q(\sigma_w^2).$$

As laid out in Subsection 2.6.1, the locality property of MFVB means that calculations concerning a particular parameter can be confined to “nearby” parameters through graph-theoretic representations of Bayesian hierarchical models. Therefore, some of the  $q$ -densities for this measurement error model are remained the same as in model (2.11). Nonetheless, we show that the optimal  $q$ -densities admit the following forms under the production restriction (4.5) in Appendix 4.A:

$$q^*(\mathbf{x}_{\text{unobs}}) \text{ is the } N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}, \boldsymbol{\sigma}_{q(\mathbf{x}_{\text{unobs},i)}}^2 \mathbf{I}_{n_{\text{unobs}}}) \text{ density function,} \quad (4.6)$$

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function,}$$

$$q^*(\mu_x) \text{ is the } N(\mu_{q(\mu_x)}, \sigma_{q(\mu_x)}^2) \text{ density function}$$

$$q^*(\sigma_x^2) \text{ is the Inverse-Gamma } \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_x^2)} \right) \text{ density function,}$$

$$q^*(\sigma_w^2) \text{ is the Inverse-Gamma } \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_w^2)} \right) \text{ density function,}$$

$$q^*(\sigma_\varepsilon^2) \text{ is the Inverse-Gamma } \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)} \right) \text{ density function,}$$

$$q^*(a_\varepsilon) \text{ is the Inverse-Gamma}(1, B_{q(a_\varepsilon)}) \text{ density function,}$$

$$q^*(\sigma_{u\ell}^2) \text{ is the Inverse-Gamma } \left( \frac{1}{2} (q_\ell^G + 1), B_{q(\sigma_{u\ell}^2)} \right) \text{ density function,}$$

$$q^*(a_{u\ell}) \text{ is the Inverse-Gamma}(1, B_{q(a_{u\ell})}) \text{ density function,}$$

$$q^*(\boldsymbol{\Sigma}^R) \text{ is the Inverse-Wishart } \left( \nu + m + q^R + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}^R)} \right) \text{ density function,}$$

$$q^*(a_r^R) \text{ is the Inverse-Gamma } \left( \frac{1}{2} (\nu + q^R), B_{q(a_r^R)} \right) \text{ density function,}$$

where the parameters are updated according to Algorithm 11, with  $\ell = 1$ . The variational lower bound on the marginal log-likelihood at the bottom of the main loop is derived in

Appendix 4.A.1.

### 4.3 Gaussian semiparametric mixed models with missing data problems

In principle, there are three missingness mechanisms: *missing completely at random*, where the missing-data mechanism is independent of the data; *missing at random*, where the missing-data mechanism depends completely on the observed  $\mathbf{y}$  but not on the missing  $\mathbf{x}$ ; and *missing not at random*, where the missing-data mechanism depends on the unobserved  $\mathbf{x}$ . Faes *et al.* (2011) treat the full array of missingness scenarios but is limited to the simplest parametric and nonparametric models. In this section we extend their work to the Gaussian semiparametric mixed models, focussing on the most challenging missing not at random scenario, where the core tenets can be elucidated without excessive notation and algebra. The other two scenarios enjoy many of the assets of the missing not at random case and therefore are omitted in this section.

We begin with the same Bayesian framework in model (2.11) and incorporate a binary random variable  $R_{ij}$  such that

$$R_{ij} = \begin{cases} 1, & \text{if } x_{ij} \text{ is observed,} \\ 0, & \text{if } x_{ij} \text{ is missing,} \end{cases}$$

for  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ . Bayesian inference for the missing data models differs according to the dependence of distribution of  $R_{ij}$  on the observed data as defined in Section 4.2. Putting all of these components together, we now consider the following model which exploits the structure described above:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{X}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}) \\ \mathbf{x}_{\text{obs}} | \mu_x, \sigma_x^2 &\sim N(\mu_x \mathbf{1}_{n_{\text{obs}}}, \sigma_x^2 \mathbf{I}), \quad \mathbf{x}_{\text{mis}} | \mu_x, \sigma_x^2 \sim N(\mu_x \mathbf{1}_{n_{\text{mis}}}, \sigma_x^2 \mathbf{I}), \\ \mathbf{R} | \boldsymbol{\phi}, \mathbf{x} &\sim \text{Bernoulli}(\{1 + \exp(-\mathbf{X}\boldsymbol{\phi})\}^{-1}), \quad \boldsymbol{\phi} \sim N(\mathbf{0}, \sigma_\phi^2 \mathbf{I}) \\ \mathbf{u} | \sigma_u^2, \boldsymbol{\Sigma}^{\text{R}} &\sim N\left(\mathbf{0}, \begin{bmatrix} \sigma_u^2 \mathbf{I}_{q^{\text{G}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma}^{\text{R}} \end{bmatrix}\right), \quad \mu_x | \sigma_x^2 \sim N(0, \sigma_x^2), \\ \sigma_x^2 | a_x &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_x), \quad a_x \sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_x^2), \\ \boldsymbol{\Sigma}^{\text{R}} | a_1^{\text{R}}, \dots, a_{q^{\text{R}}}^{\text{R}} &\sim \text{Inverse-Wishart}\left(\nu + q^{\text{R}} - 1, 2\nu \text{diag}(1/a_1^{\text{R}}, \dots, 1/a_{q^{\text{R}}}^{\text{R}})\right), \quad (4.7) \\ a_r^{\text{R}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_r^2), \quad 1 \leq r \leq q^{\text{R}}, \\ \sigma_u^2 | a_u &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_u), \quad a_u \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_u^2), \\ \sigma_w^2 | a_w &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_w), \quad a_w \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_w^2), \end{aligned}$$



### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mu_{q(1/\sigma_x^2)} > 0$ ,  $\mu_{q(1/a_x)} > 0$ ,  $\mu_{q(1/\sigma_w^2)} > 0$ ,  $\mu_{q(1/a_w)} > 0$ ,  $\mu_{q(\mu_x)} > 0$ ,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mu_{q(1/\sigma_u^2)} > 0$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 > 0$ ,  $1 \leq i \leq m$ ,  $\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})}$ , a  $n_i^{\text{unobs}} \times 1$  vector of positive entries,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , a  $(p + q^R m + q^G) \times 1$  vector and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ , a  $(p + m q^R + q^G) \times (p + m q^R + q^G)$  positive definite matrix.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q((\Omega)^{-1})} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)} \mathbf{I}_{q^G} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\Sigma^R)^{-1})} \end{bmatrix}$$

**Update Univariate Normal  $q^*(\mathbf{x}_{\text{unobs},i})$  parameters:**

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

If ( $n_i^{\text{unobs}} > 0$ ) then

$$\begin{aligned} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 &\leftarrow 1 / \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \boldsymbol{\mu}_{q(\beta_x)}^2 + \boldsymbol{\Sigma}_{q(\beta_x)} + 2(\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{1i}^R)}) \right\} \right. \\ &\quad \left. + \boldsymbol{\mu}_{q(u_{1i}^R)}^2 + \boldsymbol{\Sigma}_{q(u_{1i}^R)} + \mu_{q(1/\sigma_x^2)} + \mu_{q(1/\sigma_w^2)} \right] \\ \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})} &\leftarrow \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \mathbf{I}_{n_i^{\text{unobs}}} \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\mu}_{q(u_{1i}^R)}) \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} \right. \right. \\ &\quad - (\boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\Lambda}_{q(\beta_0, \beta_x)} + \boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_0, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{unobs}}} \\ &\quad - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\beta_s)} + \boldsymbol{\Lambda}_{q(\beta_x, \beta_s)} + \boldsymbol{\mu}_{q(\beta_s)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_s, u_{1i}^R)}) \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} \\ &\quad - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\beta_x, u_k^G)}) z_{q^G}(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) \\ &\quad - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(u_{1i}^R)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(u_{1i}^R, u_k^G)}) z_{q^G}(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) \\ &\quad \left. \left. - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{0i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{0i}^R)} + \boldsymbol{\mu}_{q(u_{0i}^R)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(u_{0i}^R, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{unobs}}} \right\} \right. \\ &\quad \left. + \mu_{q(1/\sigma_x^2)} \mu_{q(\mu_x)} \mathbf{1}_{n_i^{\text{unobs}}} + \mu_{q(1/\sigma_w^2)} \mathbf{w}_{\mathbf{x}_{\text{unobs},i}} \right] \end{aligned}$$

**Update  $E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C})$  and  $E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C}^\top \mathbf{C})$  according to Subsection 4.2.1.**

**Update multivariate normal  $q^*(\boldsymbol{\beta}, \mathbf{u})$  parameters:**

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \{ \mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C}^\top \mathbf{C}) + \mathbf{M}_{q((\Omega)^{-1})} \}^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \{ \mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C})^\top \mathbf{y} \}$$

**Update univariate normal  $q^*(\mu_x)$  parameters:**

$$\sigma_{q(\mu_x)}^2 \leftarrow 1 / (\sum_{i=1}^m n_i \mu_{q(1/\sigma_x^2)} + 1/\sigma_{\mu_x}^2)$$

$$\mu_{q(\mu_x)} \leftarrow \sigma_{q(\mu_x)}^2 \mu_{q(1/\sigma_x^2)} (\mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})})$$

### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

**Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  parameters:**

$$\begin{aligned}
B_{q(\sigma_\varepsilon^2)} &\leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left[ \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}_{x_{\text{obs}}}^\top \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \right. \\
&\quad + \|\mathbf{y}_{x_{\text{unobs}}}\|^2 - 2 \mathbf{y}_{x_{\text{unobs}}}^\top E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C}_{x_{\text{unobs}}}) \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \\
&\quad \left. + \text{tr} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{C}_{x_{\text{unobs}}}^\top \mathbf{C}_{x_{\text{unobs}}}) (\boldsymbol{\mu}_{q(\beta, \mathbf{u})} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}^\top + \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \right\} \right] \\
\mu_{q(1/\sigma_\varepsilon^2)} &\leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2})
\end{aligned}$$

**Update inverse-Gamma  $q^*(\sigma_x^2)$  parameters:**

$$\begin{aligned}
B_{q(\sigma_x^2)} &\leftarrow \mu_{q(1/a_x)} + \frac{1}{2} \left( \|\mathbf{x}_{\text{obs}} - \mu_{q(\mu_x)} \mathbf{1}\|^2 + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})} - \mu_{q(\mu_x)} \mathbf{1}\|^2 \right. \\
&\quad \left. + \sum_{i=1}^m n_i \sigma_{q(\mu_x)}^2 + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs}, i})}^2 \right) \\
\mu_{q(1/\sigma_x^2)} &\leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_x^2)} \quad ; \quad \mu_{q(1/a_x)} \leftarrow 1 / (\mu_{q(1/\sigma_x^2)} + A_x^{-2})
\end{aligned}$$

**Update inverse-Gamma  $q^*(\sigma_w^2)$  parameters:**

$$\begin{aligned}
B_{q(\sigma_w^2)} &\leftarrow \mu_{q(1/a_w)} + \frac{1}{2} \left( \|\mathbf{w}_{x_{\text{obs}}} - \mathbf{x}_{\text{obs}}\|^2 + \|\mathbf{w}_{x_{\text{unobs}}} - \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}\|^2 \right. \\
&\quad \left. + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs}, i})}^2 \right) \\
\mu_{q(1/\sigma_w^2)} &\leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_w^2)} \quad ; \quad \mu_{q(1/a_w)} \leftarrow 1 / (\mu_{q(1/\sigma_w^2)} + A_w^{-2})
\end{aligned}$$

**Update inverse-Gamma  $q^*(a_r^{\text{R}})$  and inverse-Gamma  $q^*(\boldsymbol{\Sigma}^{\text{R}})$  parameters:**

For  $r = 1, \dots, q^{\text{R}}$ :

$$\begin{aligned}
B_{q(a_r^{\text{R}})} &\leftarrow \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})})_{rr} + A_{\text{R}r}^{-2} \quad ; \quad \mu_{q(1/a_r^{\text{R}})} \leftarrow \frac{1}{2} (\nu + q^{\text{R}}) / B_{q(a_r^{\text{R}})} \\
\mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})} &\leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})} \boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^{\text{R}})}) + 2 \nu \text{diag}(\mu_{q(1/a_1^{\text{R}})}, \dots, \mu_{q(1/a_{q^{\text{R}}}^{\text{R}})}) \\
\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} &\leftarrow (\nu + m + q^{\text{R}} - 1) \mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})}^{-1}
\end{aligned}$$

**Update inverse-Gamma  $q^*(a_u)$  and inverse-Gamma  $q^*(\sigma_u^2)$  parameters:**

$$\begin{aligned}
\mu_{q(1/a_u)} &\leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \\
\mu_{q(1/\sigma_u^2)} &\leftarrow \frac{q^{\text{G}} + 1}{2 \mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{G}})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{G}})})}
\end{aligned}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**

---

Algorithm 11: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Gaussian response subject to classical measurement errors.

### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

$$\sigma_\varepsilon^2 | a_\varepsilon \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_\varepsilon\right), \quad a_\varepsilon \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_\varepsilon^2\right).$$

Again, whilst model (4.7) extends model (2.11) by including the missingness mechanism, it follows the same way to be represented as a mixed model as described in Section 2.3. The prior specification of model (4.1) can be found in Section 2.4. The DAG in Figure 4.2 shows the interplay between the regression parameters and missing data mechanism parameters.

In Section 2.7, we establish that Bayesian models with logistic regression components benefit from the introduction of the variational parameter vector  $\boldsymbol{\xi}$  via Jaakkola and Jordan's approximation (Jaakkola and Jordan, 1997). As will become clear in the Appendix 4.B that the MFVB approximation become completely algebraic (without the need of numerical quadrature or Monte Carlo methods) if the variational parameter is introduced into the model.

#### 4.3.1 Notation

The following notation is useful in the upcoming sections. Let  $n_i^{\text{obs}}$  denote the number of observed  $x_{ij}$ s in the  $i$ th group and  $n_i^{\text{unobs}}$  denote the number of missing  $x_{ij}$ s in the  $i$ th group. Let  $\mathbf{x}_{\text{obs},i}$  be the  $n_i^{\text{obs}} \times 1$  vector containing the observed  $x_{ij}$ s and  $\mathbf{x}_{\text{mis},i}$  be the  $n_i^{\text{mis}} \times 1$  containing the missing  $x_{ij}$ s. We reorder the data so that the observed data is first for each group. Hence, the full vector of predictor is

$$\mathbf{x} \equiv [\mathbf{x}_{\text{obs},1} \quad \mathbf{x}_{\text{mis},1} \quad \cdots \quad \mathbf{x}_{\text{obs},m} \quad \mathbf{x}_{\text{mis},m}]^\top.$$

In addition, let  $x_{\text{obs},ij}$  be the value of the predictor corresponding to the  $j$ th entry of  $\mathbf{x}_{\text{obs},i}$  and  $x_{\text{mis},ij'}$  be the value of the predictor corresponding to the  $j'$ th entry of  $\mathbf{x}_{\text{mis},i}$  for  $1 \leq i \leq m$ ,  $1 \leq j \leq n_i^{\text{obs}}$  and  $1 \leq j' \leq n_i^{\text{mis}}$ .

Define the following vectors and matrices

$$\begin{aligned} \mathbf{x}_i &= \begin{bmatrix} \mathbf{x}_{\text{obs},i} \\ \mathbf{x}_{\text{unobs},i} \end{bmatrix}, \quad \mathbf{x}_{\text{obs},i} = \begin{bmatrix} x_{\text{obs},i1} \\ \vdots \\ x_{\text{obs},in_i^{\text{obs}}} \end{bmatrix}, \quad \mathbf{x}_{\text{mis},i} = \begin{bmatrix} x_{\text{mis},i1} \\ \vdots \\ x_{\text{mis},in_i^{\text{mis}}} \end{bmatrix}, \\ \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} &= \text{blockdiag} \left( [\mathbf{1} \quad \mathbf{x}_{\text{obs},i}] \right)_{1 \leq i \leq m}, \quad \mathbf{Z}_{x_{\text{mis}}}^{\text{R}} = \text{blockdiag} \left( [\mathbf{1} \quad \mathbf{x}_{\text{mis},i}] \right)_{1 \leq i \leq m}, \\ \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} &= \begin{bmatrix} z_1(\mathbf{s}_{x_{\text{obs},1}}) & \cdots & z_q(\mathbf{s}_{x_{\text{obs},1}}) \\ \vdots & & \vdots \\ z_1(\mathbf{s}_{x_{\text{obs},m}}) & \cdots & z_q(\mathbf{s}_{x_{\text{obs},m}}) \end{bmatrix}, \quad \mathbf{Z}_{x_{\text{mis}}}^{\text{G}} = \begin{bmatrix} z_1(\mathbf{s}_{x_{\text{mis},1}}) & \cdots & z_q(\mathbf{s}_{x_{\text{mis},1}}) \\ \vdots & & \vdots \\ z_1(\mathbf{s}_{x_{\text{mis},m}}) & \cdots & z_q(\mathbf{s}_{x_{\text{mis},m}}) \end{bmatrix}. \end{aligned}$$

### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

For conciseness of forthcoming expressions we adopt the following notations:

$$\begin{aligned} \mathbf{C} &= \begin{bmatrix} \mathbf{C}_{x_{\text{obs}}} \\ \mathbf{C}_{x_{\text{mis}}} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} \\ \mathbf{1} & \mathbf{x}_{\text{mis}} & \mathbf{s}_{x_{\text{mis}}} & \mathbf{Z}_{x_{\text{mis}}}^{\text{G}} & \mathbf{Z}_{x_{\text{mis}}}^{\text{R}} \end{bmatrix} \\ E_q(\mathbf{C}) &= \begin{bmatrix} \mathbf{C}_{x_{\text{obs}}} \\ E_q(\mathbf{C}_{x_{\text{mis}}}) \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{G}} & \mathbf{Z}_{x_{\text{obs}}}^{\text{R}} \\ \mathbf{1} & E_q(\mathbf{x}_{\text{mis}}) & \mathbf{s}_{x_{\text{mis}}} & \mathbf{Z}_{x_{\text{mis}}}^{\text{G}} & E_q(\mathbf{Z}_{x_{\text{mis}}}^{\text{R}}) \end{bmatrix}, \end{aligned}$$

where

$$\begin{aligned} E_q(\mathbf{x}_{\text{mis}}) &= [\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},1})} \quad \cdots \quad \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},m})}]^{\text{T}} \\ \text{and } E_q(\mathbf{Z}_{x_{\text{mis}}}^{\text{R}}) &= \text{blockdiag}([\mathbf{1} \quad \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})}])_{1 \leq i \leq m}. \end{aligned}$$

The update expressions for  $E_q(\mathbf{C}^{\text{T}}\mathbf{C})$  are analogous to those in (4.3) except the subscript “unobs” is now replaced by “mis” and, in addition,

$$\begin{aligned} E_q[\mathbf{X}^{\text{T}}\text{diag}\{\lambda(\boldsymbol{\xi})\}\mathbf{X}] &= \begin{bmatrix} \mathbf{1}_{n_{\text{obs}}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{obs}}})\}\mathbf{1}_{n_{\text{obs}}} + \mathbf{1}_{n_{\text{mis}}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis}}})\}\mathbf{1}_{n_{\text{mis}}} \\ \mathbf{x}_{\text{obs}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{obs}}})\}\mathbf{x}_{\text{obs}} + \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis}}})\}\mathbf{1}_{n_{\text{mis}}} \\ \mathbf{1}_{n_{\text{obs}}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{obs}}})\}\mathbf{1}_{n_{\text{obs}}} + \mathbf{1}_{n_{\text{mis}}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis}}})\}\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})} \\ \mathbf{x}_{\text{obs}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{obs}}})\}\mathbf{x}_{\text{obs}} + E_q[\mathbf{x}_{\text{mis}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis}}})\}\mathbf{x}_{\text{mis}}] \end{bmatrix}, \quad (4.8) \end{aligned}$$

where

$$\begin{aligned} E_q[\mathbf{x}_{\text{mis}}^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis}}})\}\mathbf{x}_{\text{mis}}] &= \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})}^{\text{T}} \text{diag}\{\lambda\boldsymbol{\xi}_{x_{\text{mis}}}\}\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})} \\ &\quad + \sum_{i=1}^m \sum_{j=1}^{n_i} \lambda(\xi_{x_{\text{mis}},ij}) \sigma_{q(\mathbf{x}_{\text{mis}},i)}^2. \end{aligned}$$

#### 4.3.2 Approximate Bayesian inference via mean field variational Bayes

We now provide details on MFVB fitting and inference catered to the Gaussian response model with missing data as described in model (4.7). The essence of our MFVB approach lies in approximating the full conditional joint posterior density function of the form

$$\begin{aligned} p(\mathbf{x}_{\text{mis}}, \boldsymbol{\beta}, \boldsymbol{\phi}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mu_x, \boldsymbol{\Sigma}^{\text{R}}, \sigma_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_\phi^2 | \mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{R}) &\approx \\ q(\mathbf{x}_{\text{mis}}, \boldsymbol{\beta}, \boldsymbol{\phi}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mu_x, \boldsymbol{\Sigma}^{\text{R}}, \sigma_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_\phi^2) & \end{aligned}$$

### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

subject to the  $q$ -density product restriction

$$q(\mathbf{x}_{\text{mis}}, \boldsymbol{\beta}, \boldsymbol{\phi}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mu_x, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_\phi^2) = \quad (4.9)$$

$$q(\mathbf{x}_{\text{mis}}, \boldsymbol{\beta}, \boldsymbol{\phi}, \mathbf{u}, \mathbf{a}^{\text{R}}, \mathbf{a}_u, a_\varepsilon, \mu_x) \times q(\boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2, \sigma_\varepsilon^2, \sigma_x^2, \sigma_\phi^2).$$

Restriction (4.9) is the minimal factorisation in the MFVB approximation, deriving based on the heuristic arguments provided by Menictas and Wand (2013). One can also use the concepts of induced factorisation and moralisation introduced in Subsection 2.5.1 to establish additional product restriction

$$q(\mathbf{x}_{\text{mis}}) q(\boldsymbol{\beta}, \mathbf{u}) q(\boldsymbol{\phi}) \left\{ \prod_{r=1}^{q^{\text{R}}} q(a_r^{\text{R}}) \right\} \left\{ \prod_{\ell=1}^L q(a_{u\ell}) \right\} q(a_\varepsilon) q(\mu_x) q(\boldsymbol{\Sigma}^{\text{R}}) \left\{ \prod_{\ell=1}^L q(\sigma_{u\ell}^2) \right\} \quad (4.10)$$

$$\times q(\sigma_\varepsilon^2) q(\sigma_x^2) q(\sigma_\phi^2).$$

In Appendix 4.B, we show that the optimal  $q$ -densities admit the following forms under the production restriction (4.10):

$$\boldsymbol{\xi} \leftarrow \sqrt{\text{diagonal}\{E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X})(\boldsymbol{\Sigma}_{q(\boldsymbol{\phi})} + \boldsymbol{\mu}_{q(\boldsymbol{\phi})}\boldsymbol{\mu}_{q(\boldsymbol{\phi})}^{\text{T}})E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X})^{\text{T}}\}}, \quad (4.11)$$

$q^*(\mathbf{x}_{\text{mis}})$  is the  $N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})}, \sigma_{q(\mathbf{x}_{\text{mis}})}^2 \mathbf{I}_{n_{\text{miss}}})$  density function,

$q^*(\boldsymbol{\beta}, \mathbf{u})$  is the  $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$  density function,

$q^*(\boldsymbol{\phi})$  is the  $N(\boldsymbol{\mu}_{q(\boldsymbol{\phi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\phi})})$  density function,

$q^*(\mu_x)$  is the  $N(\mu_{q(\mu_x)}, \sigma_{q(\mu_x)}^2)$  density function

$q^*(\sigma_x^2)$  is the Inverse-Gamma  $(\frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_x^2)})$  density function,

$q^*(\sigma_\phi^2)$  is the Inverse-Gamma  $(\frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\phi^2)})$  density function,

$q^*(\sigma_\varepsilon^2)$  is the Inverse-Gamma  $(\frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_\varepsilon^2)})$  density function,

$q^*(a_\varepsilon)$  is the Inverse-Gamma  $(1, B_{q(a_\varepsilon)})$  density function,

$q^*(\sigma_{u\ell}^2)$  is the Inverse-Gamma  $(\frac{1}{2} (q_\ell^{\text{G}} + 1), B_{q(\sigma_{u\ell}^2)})$  density function,

$q^*(a_{u\ell})$  is the Inverse-Gamma  $(1, B_{q(a_{u\ell})})$  density function,

$q^*(\boldsymbol{\Sigma}^{\text{R}})$  is the Inverse-Wishart  $(\nu + m + q^{\text{R}} + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})})$  density function,

$q^*(a_r^{\text{R}})$  is the Inverse-Gamma  $(\frac{1}{2}(\nu + q^{\text{R}}), B_{q(a_r^{\text{R}})})$  density function,

where the parameters are updated according to Algorithm 12, with  $\ell = 1$ .

#### 4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mu_{q(1/\sigma_x^2)} > 0$ ,  $\mu_{q(1/a_x)} > 0$ ,  $\mu_{q(1/\sigma_w^2)} > 0$ ,  $\mu_{q(1/a_w)} > 0$ ,  $\mu_{q(1/a_r^R)} > 0$ ,  $1 \leq r \leq q^R$ ,  $\mu_{q(1/\sigma_u^2)} > 0$ ,  $\mathbf{M}_{q((\Sigma^R)^{-1})}$ , a  $q^R \times q^R$  positive definite matrix,  $\sigma_{q(\mathbf{x}_{\text{mis}})} > 0$ ,  $1 \leq i \leq m$ ,  $\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})}$ , a  $n_i^{\text{mis}} \times 1$  vector of positive entries,  $\boldsymbol{\mu}_{q(\phi)}$ , a  $2 \times 1$  vector of positive entries,  $\boldsymbol{\Sigma}_{q(\phi)}$ , a  $2 \times 2$  positive definite matrix,  $\boldsymbol{\mu}_{q(\beta, \mathbf{u})}$ , a  $(p + q^R m + q^G) \times 1$  vector and  $\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}$ , a  $(p + m q^R + q^G) \times (p + m q^R + q^G)$  positive definite matrix.

**Cycle through updates:**

$$\text{Define: } M_{q((\Omega)^{-1})} \equiv \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_\varepsilon^2)} \mathbf{I}_{q^G} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes M_{q((\Sigma^R)^{-1})} \end{bmatrix}$$

**Update Univariate Normal  $q^*(\mathbf{x}_{\text{mis},i})$  parameters:**

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

If ( $n_i^{\text{mis}} > 0$ ) then

$$\begin{aligned} \sigma_{q(\mathbf{x}_{\text{mis},i})} \leftarrow & 1 / \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \boldsymbol{\mu}_{q(\beta_x)}^2 + \boldsymbol{\Sigma}_{q(\beta_x)} + 2(\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{1i}^R)}) \right. \right. \\ & \left. \left. + \boldsymbol{\mu}_{q(u_{1i}^R)}^2 + \boldsymbol{\Sigma}_{q(u_{1i}^R)} \right\} + \mu_{q(1/\sigma_x^2)} \right. \\ & \left. + 2(\boldsymbol{\mu}_{q(\phi_1)}^2 + \boldsymbol{\Sigma}_{q(\phi_1)}) \sum_{i=1}^m \lambda(\boldsymbol{\xi}_{\mathbf{x}_{\text{mis},i}}) \right] \end{aligned}$$

$$\begin{aligned} \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})} \leftarrow & \sigma_{q(\mathbf{x}_{\text{mis},i})} \mathbf{I}_{n_i^{\text{mis}}} \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\mu}_{q(u_{1i}^R)}) \mathbf{y}_{\mathbf{x}_{\text{mis},i}} \right. \right. \\ & - (\boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\Lambda}_{q(\beta_0, \beta_x)} + \boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_0, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{mis}}} \\ & - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\beta_s)} + \boldsymbol{\Lambda}_{q(\beta_x, \beta_s)} + \boldsymbol{\mu}_{q(\beta_s)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_s, u_{1i}^R)}) \mathbf{s}_{\mathbf{x}_{\text{mis},i}} \\ & - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\beta_x, u_k^G)}) z_{q^G}(\mathbf{s}_{\mathbf{x}_{\text{mis},i}}) \\ & - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(u_{1i}^R)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(u_{1i}^R, u_k^G)}) z_{q^G}(\mathbf{s}_{\mathbf{x}_{\text{mis},i}}) \\ & \left. \left. - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{0i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{0i}^R)} + \boldsymbol{\mu}_{q(u_{0i}^R)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(u_{0i}^R, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{mis}}} \right\} \right. \\ & \left. + \mu_{q(1/\sigma_x^2)} \mu_{q(\mu_x)} \mathbf{1}_{n_i^{\text{mis}}} + 2(\boldsymbol{\mu}_{q(\phi_1)}^2 + \boldsymbol{\Sigma}_{q(\phi_1)}) \lambda(\boldsymbol{\xi}_{\mathbf{x}_{\text{mis},i}}) \right] \end{aligned}$$

**Update  $E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C})$  and  $E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C}^\top \mathbf{C})$  according to Subsection 4.3.1.**

**Update multivariate normal  $q^*(\beta, \mathbf{u})$  parameters:**

$$\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \leftarrow \{ \mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C}^\top \mathbf{C}) + M_{q((\Omega)^{-1})} \}^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u})} \leftarrow \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \{ \mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C})^\top \mathbf{y} \}$$

**Update univariate normal  $q^*(\mu_x)$  parameters:**

$$\sigma_{q(\mu_x)}^2 \leftarrow 1 / (\sum_{i=1}^m n_i \mu_{q(1/\sigma_x^2)} + 1/\sigma_{\mu_x}^2)$$

$$\mu_{q(\mu_x)} \leftarrow \sigma_{q(\mu_x)}^2 \mu_{q(1/\sigma_x^2)} (\mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})})$$

4.3. GAUSSIAN SEMIPARAMETRIC MIXED MODELS WITH MISSING DATA PROBLEMS

---

**Update inverse-Gamma  $q^*(\sigma_\varepsilon^2)$  parameters:**

$$B_{q(\sigma_\varepsilon^2)} \leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left[ \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}_{x_{\text{obs}}}^\top \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \right. \\ \left. + \|\mathbf{y}_{x_{\text{mis}}}\|^2 - 2 \mathbf{y}_{x_{\text{mis}}}^\top E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C}_{x_{\text{mis}}}) \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \right. \\ \left. + \text{tr} \left\{ E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{C}_{x_{\text{mis}}}^\top \mathbf{C}_{x_{\text{mis}}}) (\boldsymbol{\mu}_{q(\beta, \mathbf{u})} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}^\top + \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \right\} \right]$$

$$\mu_{q(1/\sigma_\varepsilon^2)} \leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_\varepsilon^2)} \quad ; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1 / (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2})$$

**Update inverse-Gamma  $q^*(\sigma_x^2)$  parameters:**

$$B_{q(\sigma_x^2)} \leftarrow \mu_{q(1/a_x)} + \frac{1}{2} \left( \|\mathbf{x}_{\text{obs}} - \mu_{q(\mu_x)} \mathbf{1}\|^2 + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis}})} - \mu_{q(\mu_x)} \mathbf{1}\|^2 \right. \\ \left. + \sum_{i=1}^m n_i \sigma_{q(\mu_x)}^2 + \sum_{i=1}^m n_i^{\text{mis}} \sigma_{q(\mathbf{x}_{\text{mis}, i})}^2 \right)$$

$$\mu_{q(1/\sigma_x^2)} \leftarrow \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_x^2)} \quad ; \quad \mu_{q(1/a_x)} \leftarrow 1 / (\mu_{q(1/\sigma_x^2)} + A_x^{-2})$$

**Update multivariate normal  $q^*(\phi)$  parameters:**

$$\boldsymbol{\Sigma}_{q(\phi)} \leftarrow \left[ 2 E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X}^\top \text{diag}\{\lambda(\boldsymbol{\xi}_{\mathbf{x}_{\text{mis}, i}})\}) \mathbf{X} + \frac{1}{\sigma_\phi^2} \mathbf{I} \right]^{-1}$$

$$\boldsymbol{\mu}_{q(\phi)} \leftarrow \boldsymbol{\Sigma}_{q(\phi)} E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X})^\top (\mathbf{R} - \frac{1}{2} \mathbf{1})$$

$$\boldsymbol{\xi} \leftarrow \sqrt{\text{diagonal}\{E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X})(\boldsymbol{\Sigma}_{q(\phi)} + \boldsymbol{\mu}_{q(\phi)} \boldsymbol{\mu}_{q(\phi)}^\top) E_{q(\mathbf{x}_{\text{mis}})}(\mathbf{X})^\top\}}$$

**Update inverse-Gamma  $q^*(a_r^{\text{R}})$  and inverse-Gamma  $q^*(\boldsymbol{\Sigma}^{\text{R}})$  parameters:**

For  $r = 1, \dots, q^{\text{R}}$ :

$$B_{q(a_r^{\text{R}})} \leftarrow \nu (\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})})_{rr} + A_{\text{R}r}^{-2} \quad ; \quad \mu_{q(1/a_r^{\text{R}})} \leftarrow \frac{1}{2} (\nu + q^{\text{R}}) / B_{q(a_r^{\text{R}})}$$

$$\mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})} \leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})} \boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})}^\top + \boldsymbol{\Sigma}_{q(\mathbf{u}_i^{\text{R}})}) + 2 \nu \text{diag}(\mu_{q(1/a_1^{\text{R}})}, \dots, \mu_{q(1/a_{q^{\text{R}}}^{\text{R}})})$$

$$\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \leftarrow (\nu + m + q^{\text{R}} - 1) \mathbf{B}_{q(\boldsymbol{\Sigma}^{\text{R}})}^{-1}$$

**Update inverse-Gamma  $q^*(a_u)$  and inverse-Gamma  $q^*(\sigma_u^2)$  parameters:**

$$\mu_{q(1/a_u)} \leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^{-2})$$

$$\mu_{q(1/\sigma_u^2)} \leftarrow \frac{q^{\text{G}} + 1}{2 \mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{G}})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{G}})})}$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

---

Algorithm 12: Mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Gaussian response subject to missing data problems.

## 4.4 Numerical evaluation

We conducted a series of comprehensive simulation studies to assess the performance of Algorithms 11 and 12 in terms of inferential accuracy, credible interval coverage and computational speed. We generated 30 datasets according to the following simulation settings. All hyperparameters are set to 10 000 and lastly we use  $q^G = 25$  interior knots when using the O’Sullivan splines with spacing as described in Wand and Ormerod (2008).

### Classical measurement error

For the measurement error model, we consider the following simulation setting:

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R)x_{ij} + f(s_{ij}), \sigma_\varepsilon^2),$$

where

$$f(s) = -\sin(2\pi s), \quad x_{ij} \stackrel{\text{ind.}}{\sim} N(\mu_x, \sigma_x^2); \quad s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1),$$

$$w_{ij} \stackrel{\text{ind.}}{\sim} N(x_{ij}, \sigma_w^2), \quad [u_{0i}^R \quad u_{1i}^R]^\top \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \Sigma^R).$$

$$\text{True values:} \quad \beta_0 = 0.58, \quad \beta_x = 1.89, \quad \sigma_\varepsilon^2 = 0.04, \quad \Sigma^R = \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix}, \quad \mu_x = 0.1,$$

where  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ . The proportion of observed  $x_{ij}$ s is varied with  $p_{\text{obs}} \in \{0.15, 0.25\}$ . The number of groups  $m$  is 50 and the within-group sample size  $n_i$  ranged between 40 and 50.

A commonly used scale-free measure of measurement error is the so-called *reliability ratio* (RR) (e.g. Pham *et al.*, 2013), defined by

$$RR = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_w^2}.$$

In this simulation setting we use the values  $RR \in \{0.9, 0.8, 0.7\}$  where  $RR = 0.9$  corresponds to a small amount of measurement error and  $RR = 0.7$  corresponds to a high amount of measurement error. In our examples We set  $\sigma_x^2 = 0.1$  and subsequently use the reliability ratio to determine  $\sigma_w^2$ .

### Missing not at random

For the missing data model, we consider the following simulation setting:

$$y_{ij} | \beta_0, \beta_x, u_{0i}^R, u_{1i}^R, \sigma_\varepsilon^2 \stackrel{\text{ind.}}{\sim} N(\beta_0 + u_{0i}^R + (\beta_x + u_{1i}^R)x_{ij} + f(s_{ij}), \sigma_\varepsilon^2),$$

where

$$f(s) = -\sin(2\pi s), \quad x_{ij} \stackrel{\text{ind.}}{\sim} N(\mu_x, \sigma_x^2),$$

$$s_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1); \quad [u_{0i}^R \quad u_{1i}^R]^\top \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \Sigma^R).$$



$$\begin{aligned} \text{True values: } \quad \beta_0 = 0.58, \quad \beta_x = 1.89, \quad \sigma_\varepsilon^2 = 0.04, \quad \Sigma^R &= \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix}, \\ \mu_x = 0.1 \quad \text{and} \quad \sigma_\varepsilon^2 \in \{0.05, 0.2, 0.8\}, \end{aligned}$$

where  $1 \leq i \leq m$  and  $1 \leq j \leq n_i$ . The number of groups  $m$  is 50 and the within-group sample size  $n_i$  ranged between 40 and 50. In this simulation setting the missingness is controlled by the two pairs of logistic coefficients:

$$\begin{aligned} (\phi_0, \phi_1) &= (2.95, -2.95) \quad \text{and} \\ (\phi_0, \phi_1) &= (0.85, -1.05), \end{aligned}$$

where the probability of missingness increases as a function of the predictor.

#### 4.4.1 Assessment of accuracy

We now focus on an empirical assessment of the accuracy of MFVB algorithms presented in Sections 4.2 and 4.3. For each simulation setting, we fitted each replicated dataset using both the MFVB and MCMC. The MFVB fits were obtained using the respective Algorithms 11 and 12 with the iterations terminated when the relative increase in  $\log p(\mathbf{y}; q)$  fell below  $10^{-7}$ . The MCMC samples were obtained using `Stan` (Stan Development Team, 2015) with R (R Development Core Team, 2015) interfacing via the `Rstan` package (Stan Development Team, 2015). In each case, MCMC samples of size 10 000 were generated. The first 5000 values of each sample were discarded as burn-in and the remaining 5000 values were thinned by a factor of 5.

As in Chapter 2 we use a measure of accuracy based on the  $L_1$  or integrated absolute error. Recall that the accuracy score, expressed as a percentage, for a particular generic parameter  $\theta$  is

$$\text{accuracy}(q^*(\theta)) = 100 \left( 1 - \frac{1}{2} \int_{-\infty}^{\infty} |q^*(\theta) - p_{\text{MCMC}}(\theta|\mathbf{y})| d\theta \right) \%.$$

Computation of  $\text{accuracy}(q^*(\theta))$  is challenging, since it depends on the posterior  $p(\theta|\mathbf{y})$  that we are trying to avoid by using MFVB methods. Instead, an approximation is made by replacing  $p(\theta|\mathbf{y})$  with a kernel density estimate  $\hat{p}(\theta|\mathbf{y})$  based on a large number of MCMC samples via the R package `KernSmooth` (Wand and Ripley, 2006).

Figure 4.3 displays the side-by-side boxplots of the accuracy scores for the key model parameters and  $f(Q_k)$ ,  $1 \leq k \leq 4$ , where  $Q_k$  are the  $k$ th sample quintiles of the  $s_{ij}$ s. For all values of  $RR$  and  $p_{\text{obs}}$ , all parameters except  $\sigma_\varepsilon^2$  and  $\sigma_w^2$  exhibit high accuracy, with almost all accuracy levels above 95%. The accuracy is poor for  $\sigma_\varepsilon^2$  and  $\sigma_w^2$ , with

performance degrading for smaller values of  $RR$  and  $p_{\text{obs}}$ .

Figure 4.4 displays the side-by-side boxplots of the accuracy scores for the key parameters and  $f(Q_k)$ ,  $1 \leq k \leq 4$ , where  $Q_k$  are the  $k$ th sample quintiles of the  $s_{ij}$ s. The parameters corresponding to the regression part of the model ( $\beta_x, \sigma_\varepsilon^2$ ) and the sample quintiles of the  $s_{ij}$ s show high accuracy across all simulation settings. However, MFVB approximations are generally poor for the missingness mechanism parameters  $\phi_0$  and  $\phi_1$  and the accuracy deteriorates as the amount of missing data gets larger or when the data becomes more noisy.

In addition, Figure 4.5 indicates that the MFVB and MCMC penalised splines fits are virtually identical, and when overlaid, are indistinguishable from one another.

#### 4.4.2 Assessment of coverage

Here we also assess the credible interval coverage against the actual coverage from the MFVB  $q$ -densities. Table 4.1 shows the percentages of true parameter coverage for the approximate 95% credible intervals formed from the MFVB densities with 0.025 probability mass in each tail, based on 1000 MFVB runs.

For the measurement error models the coverage is generally very good and does not fall below 87% for most parameters and unobservable  $x_{ij}$ s. The coverage for the variance parameters  $\sigma_x^2$  and  $\sigma_w^2$  is generally quite poor, but the coverage is better when  $p_{\text{obs}} = 0.25$ . For the missing data models the coverage is also very good overall except for the missing data mechanism parameters.

#### 4.4.3 Assessment of speed

While running the simulation studies described in Section 4.4 we kept track of the time taken for each model to be fitted and the results are summarised in Table 4.2. We used the Mac OS X laptop with a 2.6 GHz Intel Core i5 processor and 8 GBytes of random access memory. The results indicate that the MFVB approach appears to be between hundreds and thousands times faster than MCMC. However considerations such as computing environment and convergence criterion need to be taken into account to allow a fairer speed comparison.

### 4.5 Real-time mean field variational Bayes algorithms

With the ongoing advancements in the fields of electronics, computer and engineering sciences, data are not only growing rapidly in volume, but speed as well. Inevitably, this necessitates the design of statistical tools to gain insights into these high volume and velocity data, and support decision making is of key interest.

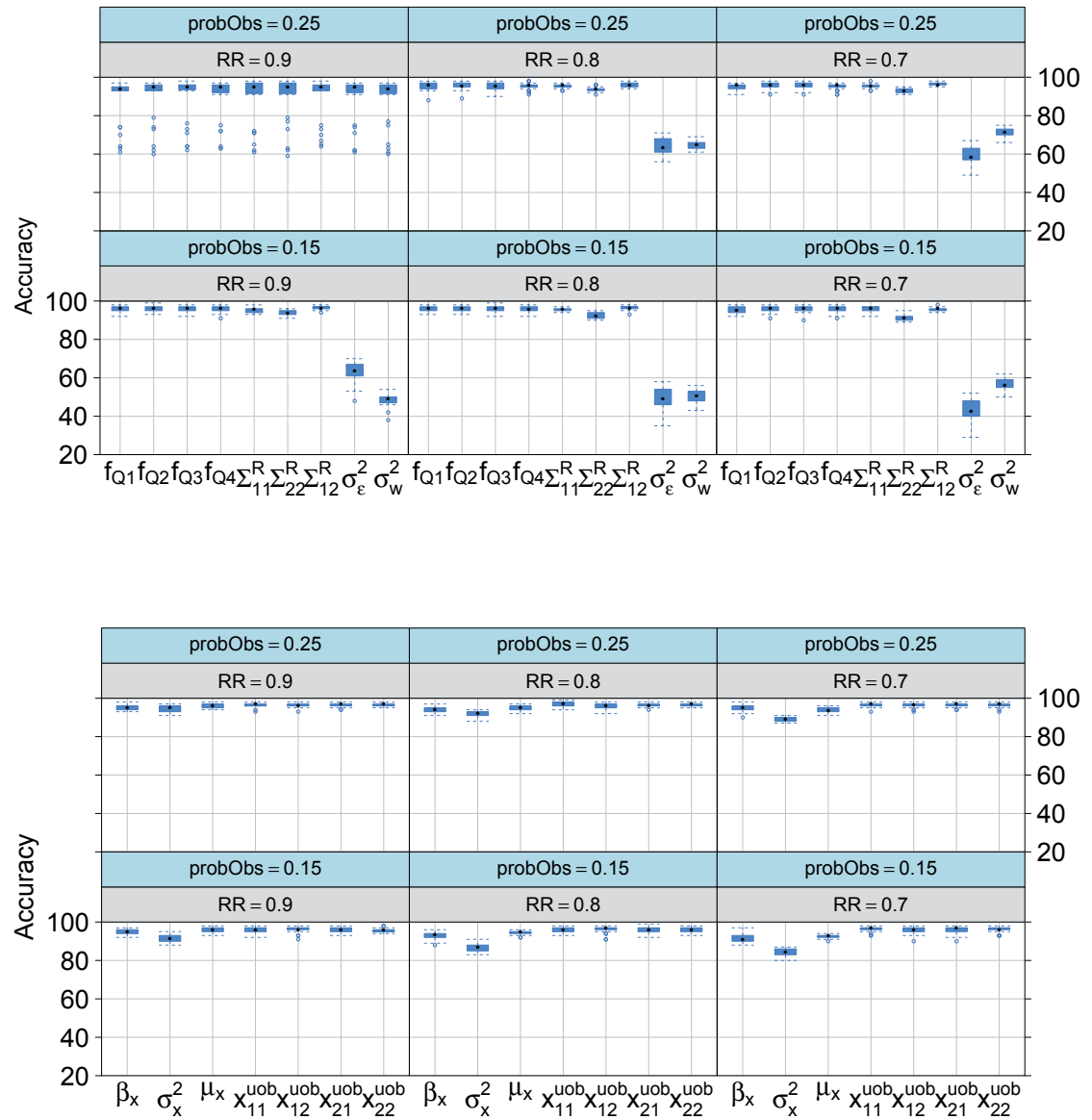


Figure 4.3: Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed model with Gaussian response subject to classical measurement errors.

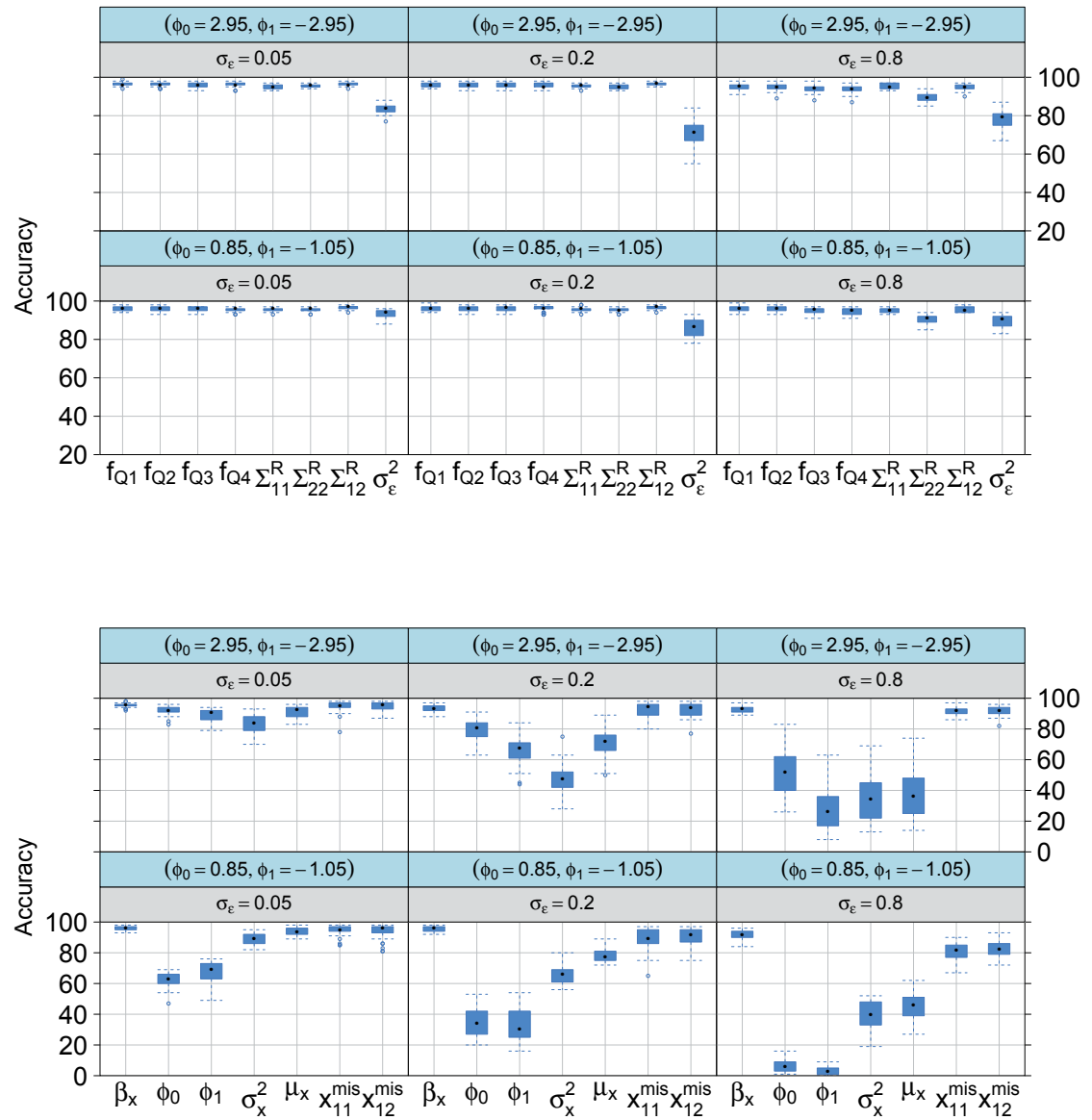


Figure 4.4: Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs for the two-level Bayesian semiparametric mixed model with Gaussian response subject to missingness.



Figure 4.5: Approximate posterior means of the penalised regression functions and corresponding pointwise 95% credible sets obtained via MFVB approximation (orange curves) and MCMC (blue curves) for a single replication of the simulation study described in the text. The red curves represent the true mean functions. The grey circles represent the observed data and the sky blue circles represent the unobserved or missing data.

#### 4.5. REAL-TIME MEAN FIELD VARIATIONAL BAYES ALGORITHMS

RR	ME $p_{\text{obs}} = 0.25$			ME $p_{\text{obs}} = 0.15$			$\sigma_\varepsilon$	MNAR low miss.			MNAR high miss.		
	0.9	0.8	0.7	0.9	0.8	0.7		0.05	0.2	0.8	0.05	0.2	0.8
$f(Q_1)$	97	97	97	97	97	97	$f(Q_1)$	96	96	94	96	96	95
$f(Q_2)$	97	97	97	97	97	97	$f(Q_2)$	96	95	95	96	95	95
$f(Q_3)$	97	97	97	97	96	97	$f(Q_3)$	96	97	97	96	97	97
$f(Q_4)$	96	96	96	95	97	96	$f(Q_4)$	96	96	96	96	96	97
$\beta_x$	97	97	97	97	97	98	$\beta_x$	98	98	97	98	98	98
$\Sigma_{11}$	94	94	94	94	94	95	$\phi_0$	92	95	67	71	39	0
$\Sigma_{12}$	95	96	95	96	96	96	$\phi_x$	92	84	4	71	32	0
$\Sigma_{22}$	96	96	96	96	96	96	$\Sigma_{11}$	95	95	96	95	95	96
$\sigma_\varepsilon^2$	81	89	65	68	47	35	$\Sigma_{12}$	93	93	93	94	94	93
$\sigma_w^2$	89	82	78	57	54	58	$\Sigma_{22}$	97	97	95	97	94	96
$\sigma_x^2$	88	87	86	86	87	83	$\sigma_\varepsilon^2$	91	82	88	94	94	95
$\mu_x$	95	94	94	93	93	91	$\sigma_x^2$	91	73	46	94	90	65
$x_{\text{unobs},11}$	91	91	94	90	92	94	$\mu_x$	91	88	95	93	94	74
$x_{\text{unobs},21}$	94	95	96	94	96	96	$x_{\text{miss},11}$	97	97	98	96	93	92
—	—	—	—	—	—	—	$x_{\text{miss},12}$	92	91	91	96	91	89

Table 4.1: Percentage coverage of true parameter values by approximate 95% credible sets based on MFVB approximate posterior density functions. A value of  $RR = 0.9$  corresponds to a small amount of measurement error and  $RR = 0.7$  corresponds to a substantial corruption of the predictor. Low missingness for the MNAR model corresponds to  $(\phi_0, \phi_1) = (2.95, -2.95)$  and high missingness to  $(\phi_0, \phi_1) = (0.85, -1.05)$ . The percentages are based on 1000 replications.

	Measurement error			Missing data		
	MCMC	MFVB	Ratio	MCMC	MFVB	Ratio
Setting 1	5424.84 (936.34)	2.13 (0.13)	2547	46277.54 (7569.84)	4.53 (0.21)	10216
Setting 2	5447.26 (864.55)	2.11 (0.28)	2582	11270.56 (491.04)	3.43 (0.26)	3286
Setting 3	5525.01 (844.71)	2.61 (0.19)	2117	2877.14 (150.42)	2.43 (0.37)	1184
Setting 4	5442.12 (673.22)	2.46 (0.48)	2212	42906.04 (2928.12)	9.36 (0.76)	4584
Setting 5	5611.41 (959.97)	3.05 (0.25)	1840	11110.27 (752.86)	8.70 (0.94)	1277
Setting 6	5548.48 (1050.81)	3.31 (0.32)	1676	2887.62 (157.70)	7.23 (0.65)	399

Table 4.2: Average (standard error) elapsed of the computing times in seconds for the MFVB and MCMC fitting of the two-level Bayesian semiparametric mixed models with Gaussian response subject to classical measurement error or missing data problem.

Up until now all proposed semiparametric regression mixed models assume that the data are processed in *batch*, that is all at the same time. The downsides to such batch processing are that the statistical analysis must wait until the entire dataset has been collected and the necessity of storing a potentially very large dataset in memory. In the *real-time* or *online* scenario, the analysis updates as each new data point (or a subset of new data points) arrive. The real-time updates use only the new data points and summary statistics from previous iterations rather than the full set of available data, thus has advantages of saving a significant amount of memory and avoiding data confidentiality.

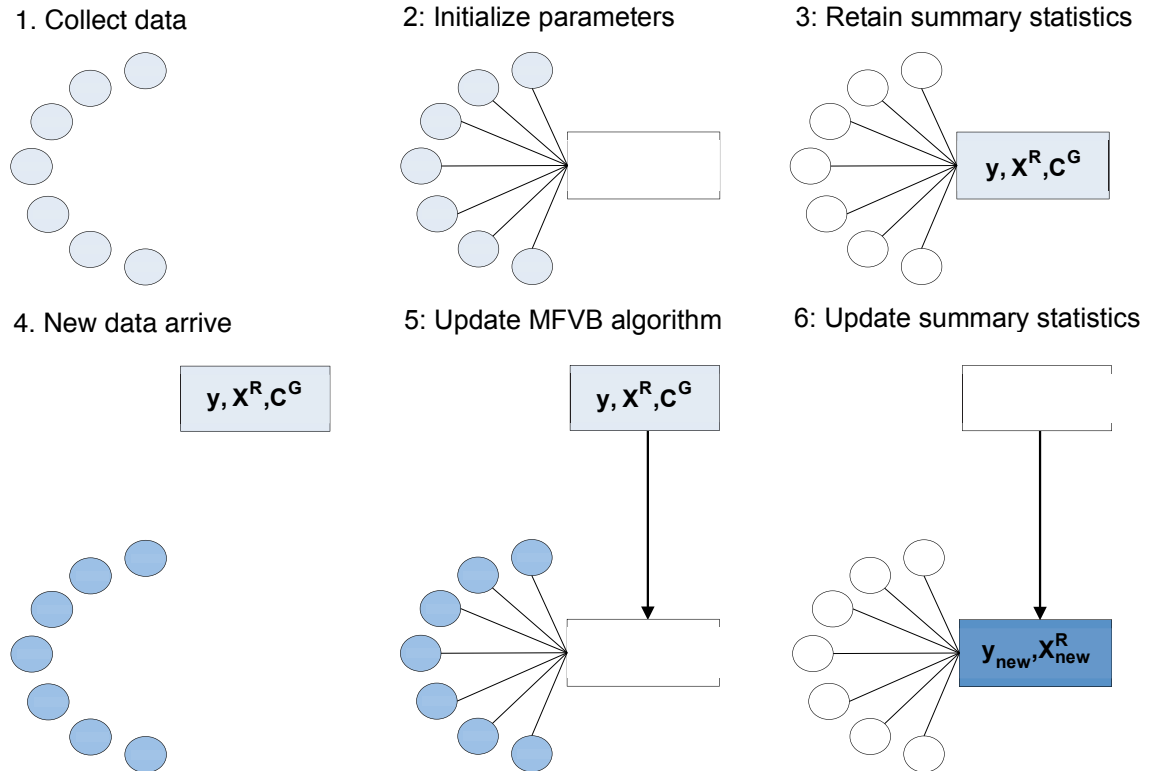


Figure 4.6: Diagrammatic description of the real-time/online variational analysis. The first phase is batch phase, where we start by running a small subset of the initial data in the batch algorithm to initialise model parameters and obtain starting values for the data sufficient statistics (these are all we keep). The second phase is the real-time/online phase where the MFVB algorithm updates as each new data point arrives. Real-time updates use only the new data and summary statistics from previous iterations rather than the full set of available data.

Figure 4.6 is a diagrammatic description of the real-time/online variational analysis.

In the machine learning literature, single-pass algorithms, one iteration per data point or per some small fixed number of data points, have recently been developed for variational inference. A recent work by Hoffman *et al.* (2010) derive a MFVB algorithm for latent Dirichlet allocation with application to topic modelling. Wang *et al.* (2011) extend their work to the hierarchical Dirichlet process. Luts *et al.* (2014) pursue a more classical definition of an “online algorithm” in that each iteration of the MFVB algorithm uses new data and past data only in the form of sufficient statistics. In this section, we utilise the approach of Luts *et al.* (2014) and develop a real-time or online version of Algorithms 3 and 4 for the longitudinal and multilevel data settings with the Gaussian or Bernoulli response. We use Algorithm 3 as the demonstrating example.

The crucial piece that allows extending Algorithm 3 towards real-time or online processing is how it uses the data: it only depends on the data through the response vector

$\mathbf{y}$ , the design matrices  $\mathbf{X}^R$  and  $\mathbf{C}^G$  and the quantities  $(\mathbf{C}^G)^\top \mathbf{y}$  and  $(\mathbf{C}^G)^\top \mathbf{C}^G$ . These quantities have simple updates when a new response vector of a new group  $\mathbf{y}_{\text{new}}$  ( $n_{\text{new}} \times 1$ ) and its corresponding  $n_{\text{new}} \times q^R$  matrix of  $\mathbf{X}_{\text{new}}^R$  and  $n_{\text{new}} \times (p + \sum_{\ell=1}^L q_\ell^G)$  matrix of  $\mathbf{C}_{\text{new}}^G$  arrive. For example, the new  $(\mathbf{C}^G)^\top \mathbf{y}$  matrix is simply

$$(\mathbf{C}^G)^\top \mathbf{y} \leftarrow (\mathbf{C}^G)^\top \mathbf{y} + (\mathbf{C}_{\text{new}}^G)^\top \mathbf{y}_{\text{new}}.$$

Continuing this way leads to Algorithm 13, the real-time or online modification of Algorithm 3. It should be noted that the starting values for the real-time procedure are determined by performing a sufficiently large batch fit. Section 2.1.1 of Luts *et al.* (2014) explains that good initialisation by the warm-up step can be important for convergence of the real-time approach. For more details regarding the batch-based tuning and convergence diagnosis process, interested readers are encouraged to refer to their paper.

The real-time Algorithm 13 differs from Algorithm 3 in that the data are processed as they arrive and the approximate posterior densities of model parameters are continually updated throughout the data collection process. Thus dynamic graphical displays as depicted in Figure 4.7 could be entertained. Figure 4.7 provides rudimentary illustration of online regression inference for the Gaussian semiparametric mixed model with data simulated according to the setting described in Section 2.10. During the warm-up phase, the number of groups is  $m_{\text{warm}} = 20$  with the total number of batch observations of 550.

The real-time or online MFVB Algorithms 13 and 14 are founded upon the same assumption as their batch counterpart. Both batch and real-time algorithms fit the Bayesian semiparametric mixed models, but the latter has the option to perform the fitting in real time for sequentially arriving data. A potential limitation with the proposed method is that the spline basis functions have to be set a priori without having the full dataset available. For example, we need to pre-specify a set of knot positions, and for many practical situations, it is plausible to specify the range of possible values beforehand, such as age. If this assumption is not reasonable, some adjustment needs to be made. In addition, we assume “fixed targets” for the model parameters, in the sense that the parameters remain more or less the same as new data arrive. Extensions to scenarios where the model parameters drift over time are beyond the scope of this thesis and would require further research. Luts (2015) explores similar extensions for the simple Gaussian linear model in the case of distributed dataset.

## 4.6 Concluding remarks

We derived efficient MFVB algorithms for inference and fitting in Gaussian semiparametric mixed models with classical measurement error and/or missing data problems. The



### Batch-based tuning and convergence diagnosis

1. Set  $m_{\text{warm}}$  to be the warm-up group size and  $m_{\text{valid}}$  to be the size of the validation period. Read in the first  $\sum_{i=1}^{m_{\text{warm}}} n_i + \sum_{i=1}^{m_{\text{valid}}} n_i$  response and predictor values.
2. Create  $\mathbf{y}_{\text{warm}} \equiv [\mathbf{y}_1 \cdots \mathbf{y}_{m_{\text{warm}}}]$ ,  $\mathbf{X}_{\text{warm}}^{\text{R}} \equiv [\mathbf{X}_1^{\text{R}} \cdots \mathbf{X}_{m_{\text{warm}}}^{\text{R}}]$  and  $\mathbf{C}_{\text{warm}}^{\text{G}} \equiv [\mathbf{C}_1^{\text{G}} \cdots \mathbf{C}_{m_{\text{warm}}}^{\text{G}}]$  consisting of the first  $m_{\text{warm}}$  group of response and predictor values.
3. Feed  $\mathbf{y}_{\text{warm}}$ ,  $\mathbf{X}_{\text{warm}}^{\text{R}}$  and  $\mathbf{C}_{\text{warm}}^{\text{G}}$  into the batch Algorithm 3 to obtain starting values for  $\mu_q(\sigma_\varepsilon^2)$ ,  $\mu_q(a_\varepsilon)$ ,  $\mathbf{M}_{q((\Sigma^{\text{R}})^{-1})}$ ,  $\mu_q(a_r^{\text{R}})$  and  $\mu_q(\sigma_{u\ell}^2)$ .
4. Set  $m \leftarrow m_{\text{warm}}$ ,  $\mathbf{n} \equiv [n_1 \cdots n_{m_{\text{warm}}}]$ ,  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow (\mathbf{C}_{\text{warm}}^{\text{G}})^{\text{T}} \mathbf{y}_{\text{warm}}$  and  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{C}^{\text{G}} \leftarrow (\mathbf{C}_{\text{warm}}^{\text{G}})^{\text{T}} \mathbf{C}_{\text{warm}}^{\text{G}}$ .
5. Run the following online MFVB Algorithm until  $m = m_{\text{warm}} + m_{\text{valid}}$ .

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_\varepsilon^2)} > 0$ ,  $\mu_{q(1/a_\varepsilon)} > 0$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $\mu_{q(1/a_u)} > 0$ ,  $\mu_{q(1/a_r^{\text{R}})} > 0$ ,  $1 \leq r \leq q^{\text{R}}$ ,  $\mathbf{M}_{q((\Sigma^{\text{R}})^{-1})}$  positive definite.

**Set:**  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow 0$  and  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{C}^{\text{G}} \leftarrow 0$ .

**Cycle:**

Read in  $\mathbf{y}_{\text{new}}$  ( $n_{\text{new}} \times 1$ ),  $\mathbf{X}_{\text{new}}^{\text{R}}$  ( $n_{\text{new}} \times q^{\text{R}}$ ) and  $\mathbf{C}_{\text{new}}^{\text{G}} \equiv [\mathbf{X}_{\text{new}} \quad \mathbf{Z}_{\text{new}}^{\text{G}}] \{(n_{\text{new}} \times (p + q^{\text{G}}))\}$

$m \leftarrow m + 1$  ;  $\mathbf{n} \leftarrow [\mathbf{n} \quad n_{\text{new}}]$  ;  $\mathbf{y} \leftarrow [\mathbf{y}^{\text{T}} \quad \mathbf{y}_{\text{new}}^{\text{T}}]^{\text{T}}$

$\mathbf{X}^{\text{R}} \leftarrow [(\mathbf{X}^{\text{R}})^{\text{T}} \quad (\mathbf{X}_{\text{new}}^{\text{R}})^{\text{T}}]^{\text{T}}$  ;  $\mathbf{C}^{\text{G}} \leftarrow [(\mathbf{C}^{\text{G}})^{\text{T}} \quad (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}}]^{\text{T}}$

$(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow (\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} + (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}} \mathbf{y}_{\text{new}}$

$(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{C}^{\text{G}} \leftarrow (\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{C}^{\text{G}} + (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}} \mathbf{C}_{\text{new}}^{\text{G}}$

Repeat update expressions for the batch MFVB Algorithm 3

6. Use convergence diagnostic graphics to assess whether the online parameters are converging to the batch parameters.
  - (a) If not converging then return to Step 1 and increase  $m_{\text{warm}}$ .
  - (b) If converging then continue running the online MFVB Algorithm until data no longer available or analysis terminated.

---

Algorithm 13: Online mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Gaussian response.

### Batch-based tuning and convergence diagnosis

1. Set  $m_{\text{warm}}$  to be the warm-up group size and  $m_{\text{valid}}$  to be the size of the validation period. Read in the first  $\sum_{i=1}^{m_{\text{warm}}} n_i + \sum_{i=1}^{m_{\text{valid}}} n_i$  response and predictor values.
2. Create  $\mathbf{y}_{\text{warm}} \equiv [\mathbf{y}_1 \cdots \mathbf{y}_{m_{\text{warm}}}]$ ,  $\boldsymbol{\xi}_{\text{warm}} \equiv [\boldsymbol{\xi}_1 \cdots \boldsymbol{\xi}_{m_{\text{warm}}}]$ ,  $\mathbf{X}_{\text{warm}}^{\text{R}} \equiv [\mathbf{X}_1^{\text{R}} \cdots \mathbf{X}_{m_{\text{warm}}}^{\text{R}}]$  and  $\mathbf{C}_{\text{warm}}^{\text{G}} \equiv [\mathbf{C}_1^{\text{G}} \cdots \mathbf{C}_{m_{\text{warm}}}^{\text{G}}]$  consisting of the first  $m_{\text{warm}}$  group of response and predictor values.
3. Feed  $\mathbf{y}_{\text{warm}}$ ,  $\mathbf{X}_{\text{warm}}^{\text{R}}$  and  $\mathbf{C}_{\text{warm}}^{\text{G}}$  into the batch Algorithm 4 to obtain starting values for  $\mathbf{G}$ ,  $\mathbf{H}$ ,  $\boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})}$ ,  $\boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})}$ ,  $\boldsymbol{\Sigma}_{q(\beta, \mathbf{u}^{\text{G}})}$ ,  $\boldsymbol{\Sigma}_{q(\mathbf{u}_i^{\text{R}})}$ ,  $\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})}$ ,  $\mu_{q(a_r^{\text{R}})}$  and  $\mu_{q(\sigma_{u\ell}^2)}$ .
4. Set  $m \leftarrow m_{\text{warm}}$ ,  $\mathbf{n} \equiv [n_1 \cdots n_{m_{\text{warm}}}]$  and  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow (\mathbf{C}_{\text{warm}}^{\text{G}})^{\text{T}} \mathbf{y}_{\text{warm}}$ .
5. Run the following online MFVB Algorithm until  $m = m_{\text{warm}} + m_{\text{valid}}$ .

**Initialise:**  $\nu = 2$ ,  $\mu_{q(1/\sigma_{u\ell}^2)} > 0$ ,  $\mu_{q(1/a_u)} > 0$ ,  $\mu_{q(1/a_r^{\text{R}})} > 0$ ,  $1 \leq r \leq q^{\text{R}}$ ,  $\mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})}$  positive definite,  $\boldsymbol{\xi}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries.

**Set:**  $(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow 0$ .

**Cycle:**

Read in  $\mathbf{y}_{\text{new}}$  ( $n_{\text{new}} \times 1$ ),  $\mathbf{X}_{\text{new}}^{\text{R}}$  ( $n_{\text{new}} \times q^{\text{R}}$ ) and  $\mathbf{C}_{\text{new}}^{\text{G}} \equiv [\mathbf{X}_{\text{new}} \quad \mathbf{Z}_{\text{new}}^{\text{G}}] \{(n_{\text{new}} \times (p + q^{\text{G}}))\}$

$\boldsymbol{\xi}_{\text{new}}^2 \leftarrow \text{diagonal}\{\mathbf{C}_{\text{new}}^{\text{G}} (\boldsymbol{\Sigma}_{q(\beta, \mathbf{u}^{\text{G}})} + \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})} \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})}^{\text{T}}) (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}}\}$

$\boldsymbol{\xi}_{\text{new}}^2 \leftarrow \boldsymbol{\xi}_{\text{new}}^2 + 2 \text{diagonal}\{\mathbf{C}_{\text{new}}^{\text{G}} (-\boldsymbol{\Sigma}_{q(\beta, \mathbf{u}^{\text{G}})} \mathbf{G}_m \mathbf{H}_m + \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})} \boldsymbol{\mu}_{q(\mathbf{u}_m^{\text{R}})}^{\text{T}}) (\mathbf{X}_{\text{new}}^{\text{R}})^{\text{T}}\}$

$\boldsymbol{\xi}_{\text{new}}^2 \leftarrow \boldsymbol{\xi}_{\text{new}}^2 + \text{diagonal}\{(\mathbf{X}_{\text{new}}^{\text{R}} (\boldsymbol{\Sigma}_{q(\mathbf{u}_m^{\text{R}})} + \boldsymbol{\mu}_{q(\mathbf{u}_m^{\text{R}})} \boldsymbol{\mu}_{q(\mathbf{u}_m^{\text{R}})}^{\text{T}}) (\mathbf{X}_{\text{new}}^{\text{R}})^{\text{T}}\}$

$m \leftarrow m + 1$  ;  $\mathbf{n} \leftarrow [\mathbf{n} \quad n_{\text{new}}]$

$\mathbf{y} \leftarrow [\mathbf{y}^{\text{T}} \quad \mathbf{y}_{\text{new}}^{\text{T}}]^{\text{T}}$  ;  $\boldsymbol{\xi} \leftarrow [\boldsymbol{\xi}^{\text{T}} \quad \boldsymbol{\xi}_{\text{new}}^{\text{T}}]^{\text{T}}$

$\mathbf{X}^{\text{R}} \leftarrow [(\mathbf{X}^{\text{R}})^{\text{T}} \quad (\mathbf{X}_{\text{new}}^{\text{R}})^{\text{T}}]^{\text{T}}$  ;  $\mathbf{C}^{\text{G}} \leftarrow [(\mathbf{C}^{\text{G}})^{\text{T}} \quad (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}}]^{\text{T}}$

$(\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} \leftarrow (\mathbf{C}^{\text{G}})^{\text{T}} \mathbf{y} + (\mathbf{C}_{\text{new}}^{\text{G}})^{\text{T}} \mathbf{y}_{\text{new}}$

Repeat update expressions for the batch MFVB Algorithm 4

6. Use convergence diagnostic graphics to assess whether the online parameters are converging to the batch parameters.
  - (a) If not converging then return to Step 1 and increase  $m_{\text{warm}}$ .
  - (b) If converging then continue running the online MFVB Algorithm until data no longer available or analysis terminated.

---

Algorithm 14: Online mean field variational Bayes algorithm for obtaining model parameters in the optimal  $q$ -density functions for the two-level Bayesian semiparametric mixed model with Bernoulli response.

4.6. CONCLUDING REMARKS

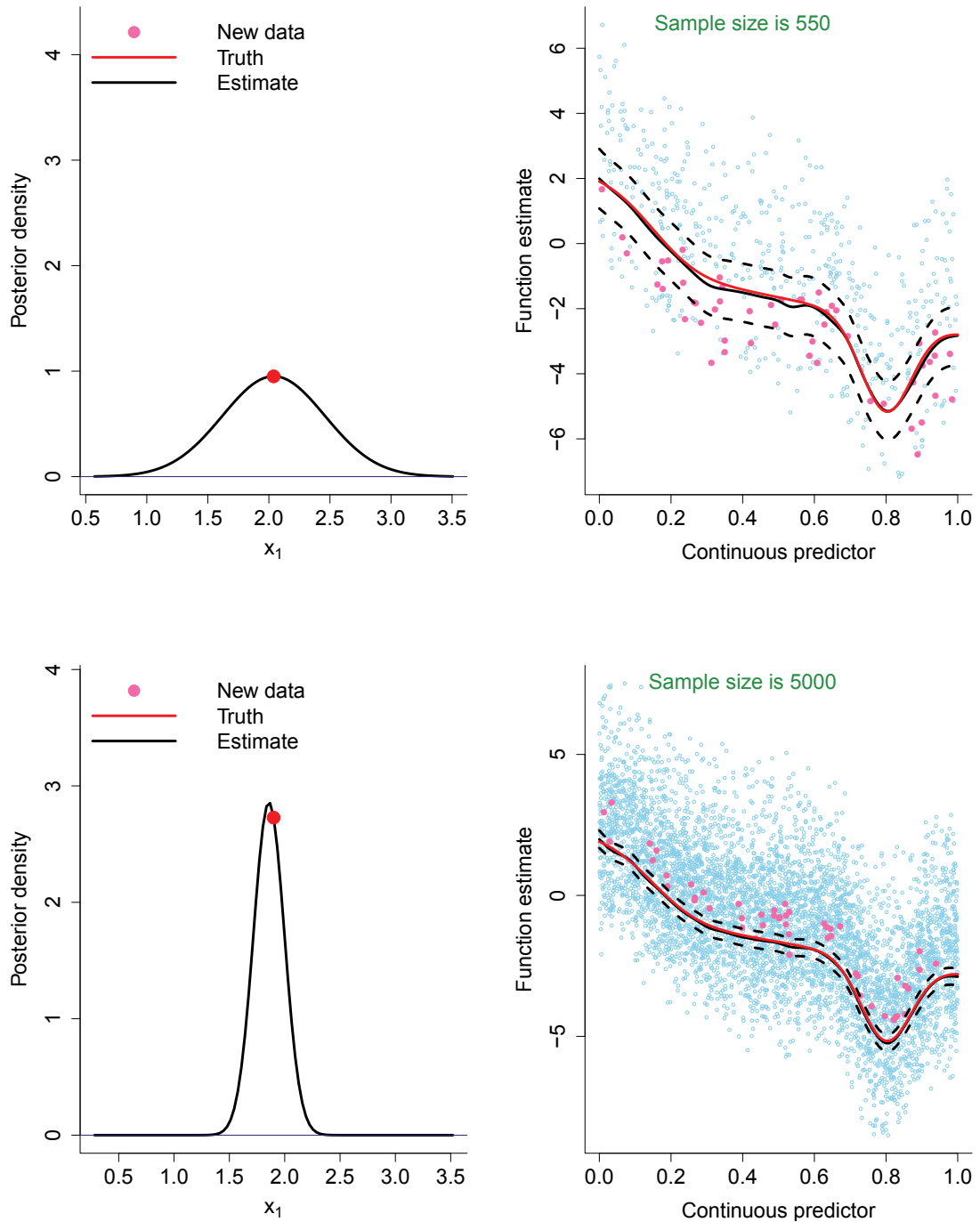


Figure 4.7: Real-time MFVB fitting of the two-level Bayesian semiparametric mixed models with Gaussian response for two different sample sizes. The sky blue circles represent the batch data and the pink circles represent the new data.

#### 4.6. CONCLUDING REMARKS

---

MFVB algorithms were shown to have moderately good to excellent accuracy for the main parameters of interest. Poor accuracy is realised for the missing data mechanism parameters, but can potentially be improved via a less stringent, but more computationally demanding mean-field restriction. In addition, the speed comparisons illustrate that MFVB achieves great improvement in efficiency over the MCMC approach for static data. Further, we showed that the advantage of MFVB approaches to approximate inference is their adaptability to real-time or online manner by combining summary statistics instead of the actual data. This aspect is particularly important in the era of high volume and velocity data.

## 4.A Optimal $q$ -densities derivation for measurement error problems

We now derive Algorithms 11 and 12 concerning MFVB fitting of the Bayesian semiparametric mixed models with data that are subject to either measurement error or missing data problems. Throughout these appendices, we use the notation “const” to denote additive constants with respect to the function argument.

We recall that

$$N \equiv \sum_{i=1}^m n_i, \quad \mathbf{v} \equiv \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \quad \mathbf{C} \equiv [\mathbf{X} \quad \mathbf{Z}]$$

$$\text{and } \boldsymbol{\Omega} \equiv \begin{bmatrix} \sigma_\beta^2 \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sigma_u^2 \mathbf{I}_{q^G} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma}^R \end{bmatrix}.$$

We first write out the full model expression

$$y_{ij} = \beta_0 + \beta_x x_{ij} + \beta_s s_{ij} + \sum_{k=1}^{q^G} u_k^G z_k(s_{ij}) + u_{0i}^R + u_{1i}^R x_{ij} + \varepsilon_{ij},$$

which can be broken into the observed and unobserved parts:

$$\begin{bmatrix} \mathbf{y}_{x_{\text{obs}}} \\ \mathbf{y}_{x_{\text{unobs}}} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{n_{\text{obs}}} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} \\ \mathbf{1}_{n_{\text{unobs}}} & \mathbf{x}_{\text{unobs}} & \mathbf{s}_{x_{\text{unobs}}} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_x \\ \beta_s \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_{x_{\text{obs}}}^G \\ \mathbf{Z}_{x_{\text{unobs}}}^G \end{bmatrix} \begin{bmatrix} u_1^G \\ \vdots \\ u_m^G \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_{x_{\text{obs}}}^R \\ \mathbf{Z}_{x_{\text{unobs}}}^R \end{bmatrix} \begin{bmatrix} \mathbf{u}_1^R \\ \vdots \\ \mathbf{u}_m^R \end{bmatrix}.$$

Expressions for  $q^*(\boldsymbol{\beta}, \mathbf{u})$ ,  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$  and  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \sim N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}),$$

where

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} = \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} E_q(\mathbf{C}^\top) \mathbf{y} \quad \text{and}$$

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} = \left( \mu_{q(1/\sigma_\varepsilon^2)} E_q(\mathbf{C}^\top \mathbf{C}) + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)} \mathbf{I}_{q^G} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})} \end{bmatrix} \right)^{-1}.$$

*Derivation:*

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

First we note that

$$\begin{aligned} p(\boldsymbol{\beta}, \mathbf{u}|\text{rest}) &\propto p(\boldsymbol{\beta}, \mathbf{u}|\mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_\varepsilon^2, \boldsymbol{\Sigma}^{\text{R}}, \boldsymbol{\sigma}_u^2) \\ &\propto p(\mathbf{y}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \sigma_\varepsilon^2) p(\boldsymbol{\beta}) p(\mathbf{u}^{\text{R}}|\boldsymbol{\Sigma}^{\text{R}}) p(\mathbf{u}^{\text{G}}|\boldsymbol{\sigma}_u^2). \end{aligned}$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(\boldsymbol{\beta}, \mathbf{u}|\text{rest}) &= -\frac{1}{2\sigma_\varepsilon^2}(\mathbf{y} - \mathbf{C}\mathbf{v})^\top(\mathbf{y} - \mathbf{C}\mathbf{v}) - \frac{1}{2}\mathbf{v}^\top\boldsymbol{\Omega}^{-1}\mathbf{v} + \text{const} \\ &= -\frac{1}{2}\left\{\mathbf{v}^\top(\sigma_\varepsilon^{-2}\mathbf{C}^\top\mathbf{C} + \boldsymbol{\Omega}^{-1})\mathbf{v} - 2\sigma_\varepsilon^{-2}\mathbf{v}^\top\mathbf{C}^\top\mathbf{y}\right\} + \text{const}, \end{aligned}$$

which then leads to the full conditional

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} \sim N\left((\sigma_\varepsilon^{-2}\mathbf{C}^\top\mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1}\sigma_\varepsilon^{-2}\mathbf{C}^\top\mathbf{y}, (\sigma_\varepsilon^{-2}\mathbf{C}^\top\mathbf{C} + \boldsymbol{\Omega}^{-1})^{-1}\right).$$

Taking expectations with respect to all parameters except  $(\boldsymbol{\beta}, \mathbf{u})$  gives

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}) &= E_q\{\log p(\boldsymbol{\beta}, \mathbf{u}|\text{rest})\} + \text{const} \\ &= -\frac{1}{2}\left\{\mathbf{v}^\top\left(\mu_{q(1/\sigma_\varepsilon^2)}E_q(\mathbf{C}^\top\mathbf{C}) + \mathbf{M}_{q(\boldsymbol{\Omega}^{-1})}\right)\mathbf{v} - 2\sigma_\varepsilon^{-2}\mathbf{v}^\top E_q(\mathbf{C}^\top)\mathbf{y}\right\} + \text{const}, \end{aligned}$$

where  $E_q(\mathbf{C}^\top)$  and  $E_q(\mathbf{C}^\top\mathbf{C})$  are given in Section 4.2.1 and

$$\mathbf{M}_{q(\boldsymbol{\Omega}^{-1})} \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)}\mathbf{I}_{q^{\text{G}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \mathbf{M}_{q((\boldsymbol{\Sigma}^{\text{R}})^{-1})} \end{bmatrix}.$$

The stated results follow via standard ‘‘completion of square’’ manipulations.

**Expression for  $q^*(\mathbf{x}_{\text{unobs},i})$ ,  $\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})}$  and  $\sigma_{q(\mathbf{x}_{\text{unobs},i})}^2$**

For  $i = 1, \dots, m$

$$q^*(\mathbf{x}_{\text{unobs},i}) \sim N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})}, \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \mathbf{I}),$$

where

$$\sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 = E_q\left(\frac{(\beta_x + u_{1i}^{\text{R}})^2}{\sigma_\varepsilon^2} + \frac{1}{\sigma_x^2} + \frac{1}{\sigma_w^2}\right)^{-1}$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

$$\begin{aligned} \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs},i}) &= \sigma_q^2(\mathbf{x}_{\text{unobs},i}) \mathbf{I}_{n_i^{\text{unobs}}} E_q \left\{ \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^R) \left( \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} - \beta_0 \mathbf{1}_{n_i^{\text{unobs}}} \right. \right. \\ &\quad \left. \left. - \beta_s \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} - \sum_{k=1}^{q^G} u_k^G z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) - u_{0i}^R \mathbf{1}_{n_i^{\text{unobs}}} \right) \right\}. \end{aligned}$$

*Derivation:*

First we note that

$$\begin{aligned} p(\mathbf{x}_{\text{unobs}}|\text{rest}) &\propto p(\mathbf{x}_{\text{unobs}}|\mathbf{x}_{\text{obs}}, \mathbf{y}, \mathbf{u}^R, \mathbf{u}^G, \boldsymbol{\beta}, \sigma_\varepsilon^2, \mathbf{w}_{\text{unobs}}, \sigma_w^2, \mu_x, \sigma_x^2) \\ &\propto p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \mathbf{u}^R, \mathbf{u}^G, \sigma_\varepsilon^2) p(\mathbf{x}_{\text{unobs}}|\mu_x, \sigma_x^2) p(\mathbf{w}_{\text{unobs}}|\mathbf{x}_{\text{unobs}}, \sigma_w^2). \end{aligned}$$

Taking logarithms on both sides where  $x$  is unobserved gives

$$\begin{aligned} \log p(\mathbf{x}_{\text{unobs}}|\text{rest}) &= \sum_{i=1}^m \left\{ -\frac{1}{2\sigma_\varepsilon^2} \left\| \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} - \left( \beta_0 \mathbf{1}_{n_i^{\text{unobs}}} + \beta_x \mathbf{x}_{\text{unobs},i} + \beta_s \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} \right. \right. \right. \\ &\quad \left. \left. + \sum_{k=1}^{q^G} u_k^G z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) + u_{0i}^R \mathbf{1}_{n_i^{\text{unobs}}} + u_{1i}^R \mathbf{x}_{\text{unobs},i} \right\|^2 \right. \\ &\quad \left. - \frac{1}{2\sigma_x^2} \left\| \mathbf{x}_{\text{unobs},i} - \mu_x \mathbf{1}_{n_i^{\text{unobs}}} \right\|^2 - \frac{1}{2\sigma_w^2} \left\| \mathbf{w}_{\mathbf{x}_{\text{unobs},i}} - \mathbf{x}_{\text{unobs},i} \right\|^2 \right\} + \text{const} \\ &= \sum_{i=1}^m \left[ -\frac{1}{2\sigma_\varepsilon^2} \left\{ (\beta_x + u_{1i}^R)^2 \left\| \mathbf{x}_{\text{unobs},i} \right\|^2 - 2(\beta_x + u_{1i}^R) \mathbf{x}_{\text{unobs},i}^\top \right. \right. \\ &\quad \left. \left. \times \left( \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} - \beta_0 \mathbf{1}_{n_i^{\text{unobs}}} - \beta_s \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} - \sum_{k=1}^{q^G} u_k^G z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) - u_{0i}^R \mathbf{1}_{n_i^{\text{unobs}}} \right) \right\} \right. \\ &\quad \left. - \frac{1}{2\sigma_x^2} \left( \left\| \mathbf{x}_{\text{unobs},i} \right\|^2 - 2 \mathbf{x}_{\text{unobs},i}^\top \mu_x \mathbf{1}_{n_i^{\text{unobs}}} \right) \right. \\ &\quad \left. - \frac{1}{2\sigma_w^2} \left( \left\| \mathbf{x}_{\text{unobs},i} \right\|^2 - 2 \mathbf{x}_{\text{unobs},i}^\top \mathbf{w}_{\mathbf{x}_{\text{unobs},i}} \right) \right] + \text{const}. \end{aligned}$$

Taking expectations with respect to all parameters except  $\mathbf{x}_{\text{unobs},i}$  gives

$$\begin{aligned} E_q \{ \log p(\mathbf{x}_{\text{unobs},i}|\text{rest}) \} &= -\frac{1}{2} \left[ \mathbf{x}_{\text{unobs},i}^\top E_q \left( \frac{(\beta_x + u_{1i}^R)^2}{\sigma_\varepsilon^2} + \frac{1}{\sigma_x^2} + \frac{1}{\sigma_w^2} \right) \mathbf{x}_{\text{unobs},i} \right. \\ &\quad - 2 \mathbf{x}_{\text{unobs},i}^\top \left\{ E_q \left( \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^R) \left( \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} - \beta_0 \mathbf{1}_{n_i^{\text{unobs}}} \right. \right. \right. \\ &\quad \left. \left. - \beta_s \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} - \sum_{k=1}^{q^G} u_k^G z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) - u_{0i}^R \right) \right\} \\ &\quad \left. + \frac{\mu_x \mathbf{1}_{n_i^{\text{unobs}}}}{\sigma_x^2} + \frac{\mathbf{w}_{\mathbf{x}_{\text{unobs},i}}}{\sigma_w^2} \right] + \text{const}, \end{aligned}$$

where

$$E_q \left( \frac{(\beta_x + u_{1i}^R)^2}{\sigma_\varepsilon^2} + \frac{1}{\sigma_x^2} + \frac{1}{\sigma_w^2} \right)$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

$$\begin{aligned}
&= \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \boldsymbol{\mu}_{q(\beta_x)}^2 + (\boldsymbol{\Sigma}_{q(\beta)})_{22} + 2(\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, \mathbf{u}_{1i}^R)}) \right\} \right. \\
&\quad \left. + \boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)}^2 + \boldsymbol{\Sigma}_{q(\mathbf{u}_{1i}^R)} + \mu_{q(1/\sigma_x^2)} + \mu_{q(1/\sigma_w^2)} \right] \\
\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})} &= \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \mathbf{I}_{n_i^{\text{unobs}}} \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)}) \mathbf{y}_{\mathbf{x}_{\text{unobs},i}} \right. \right. \\
&\quad - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\beta_0)} + \boldsymbol{\Lambda}_{q(\beta_0, \beta_x)} + \boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)} \boldsymbol{\mu}_{q(\beta_0)} + \boldsymbol{\Lambda}_{q(\mathbf{u}_{1i}^R, \beta_0)}) \mathbf{1}_{n_i^{\text{unobs}}} \\
&\quad - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\beta_s)} + \boldsymbol{\Lambda}_{q(\beta_x, \beta_s)} + \boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)} \boldsymbol{\mu}_{q(\beta_s)} + \boldsymbol{\Lambda}_{q(\mathbf{u}_{1i}^R, \beta_s)}) \mathbf{s}_{\mathbf{x}_{\text{unobs},i}} \\
&\quad - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\beta_x, u_k^G)}) z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) \\
&\quad - \sum_{k=1}^{q^G} (\boldsymbol{\mu}_{q(\mathbf{u}_{1i}^R)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\mathbf{u}_{1i}^R, u_k^G)}) z_k(\mathbf{s}_{\mathbf{x}_{\text{unobs},i}}) \\
&\quad \left. \left. - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\mathbf{u}_{0i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, \mathbf{u}_{0i}^R)} + \mu_{q(\mathbf{u}_{1i}^R)} \boldsymbol{\mu}_{q(\mathbf{u}_{0i}^R)} + \boldsymbol{\Lambda}_{q(\mathbf{u}_{0i}^R, \mathbf{u}_{1i}^R)}) \mathbf{1}_{n_i^{\text{unobs}}} \right\} \right. \\
&\quad \left. + \mu_{q(1/\sigma_x^2)} \boldsymbol{\mu}_{q(\mu_x)} \mathbf{1}_{n_i^{\text{unobs}}} + \mu_{q(1/\sigma_w^2)} \mathbf{w}_{\mathbf{x}_{\text{unobs},i}} \right].
\end{aligned}$$

**Expressions for  $q^*(\sigma_\varepsilon^2)$ ,  $B_q(\sigma_\varepsilon^2)$  and  $\mu_{q(1/\sigma_\varepsilon^2)}$**

$$q^*(\sigma_\varepsilon^2) \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_q(\sigma_\varepsilon^2) \right),$$

where

$$B_q(\sigma_\varepsilon^2) = \frac{1}{2} E_q(\|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2) + \mu_{q(1/\sigma_\varepsilon^2)} \quad \text{and} \quad \mu_{q(\sigma_\varepsilon^2)} = \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_q(\sigma_\varepsilon^2)$$

*Derivation:*

First we note that

$$\begin{aligned}
p(\sigma_\varepsilon^2 | \text{rest}) &\propto p(\sigma_\varepsilon^2 | \mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, a_\varepsilon) \\
&\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_\varepsilon^2) p(\sigma_\varepsilon^2 | a_\varepsilon).
\end{aligned}$$

Taking logarithms on both sides gives

$$\log p(\sigma_\varepsilon^2 | \text{rest}) = \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \log(\sigma_\varepsilon^2) - \left( \frac{1}{2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 + a_\varepsilon^{-1} \right) / \sigma_\varepsilon^2 + \text{const}$$

which then leads to the full conditional

$$\sigma_\varepsilon^2 | \text{rest} \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), \frac{1}{2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 + a_\varepsilon^{-1} \right).$$



#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

Taking expectations with respect to all parameters except  $\sigma_\varepsilon^2$  gives

$$\begin{aligned}\log q^*(\sigma_\varepsilon^2) &= E_q \{ \log p(\sigma_\varepsilon^2 | \text{rest}) \} + \text{const} \\ &= \left\{ -\frac{1}{2} (\sum_{i=1}^m n_i + 1) - 1 \right\} \log(\sigma_\varepsilon^2) - \left[ \frac{1}{2} \{ E_q(\|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2) \} + \mu_{q(1/a_\varepsilon)} \right] / \sigma_\varepsilon^2 + \text{const}.\end{aligned}$$

The form of  $q^*(\sigma_\varepsilon^2)$  and  $B_{q(\sigma_\varepsilon^2)}$  follows from Result 1.11.6 where

$$\begin{aligned}E_q(\|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2) &= \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{C}\mathbf{x}_{\text{obs}}\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\mathbf{C}_{x_{\text{obs}}}^\top \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \\ &\quad + \|\mathbf{y}_{x_{\text{unobs}}}\|^2 - 2\mathbf{y}_{x_{\text{unobs}}}^\top E_q(\mathbf{C}_{x_{\text{unobs}}})\boldsymbol{\mu}_{q(\mathbf{v})} \\ &\quad + \text{tr}\{E_q(\mathbf{C}_{x_{\text{unobs}}}^\top \mathbf{C}_{x_{\text{unobs}}})(\boldsymbol{\mu}_{q(\mathbf{v})}\boldsymbol{\mu}_{q(\mathbf{v})}^\top + \boldsymbol{\Sigma}_{q(\mathbf{v})})\}.\end{aligned}$$

**Expressions for  $q^*(a_\varepsilon)$ ,  $B_{q(a_\varepsilon)}$  and  $\mu_{q(1/a_\varepsilon)}$**

$$q^*(a_\varepsilon) \sim \text{Inverse-Gamma}(1, B_{q(a_\varepsilon)}),$$

where

$$B_{q(a_\varepsilon)} = \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2} \quad \text{and} \quad \mu_{q(1/a_\varepsilon)} = 1 / \{ \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2} \}.$$

*Derivation:*

First we note that

$$p(a_\varepsilon | \text{rest}) \propto p(a_\varepsilon | \sigma_\varepsilon^2) = p(\sigma_\varepsilon^2 | a_\varepsilon) p(a_\varepsilon).$$

Taking logarithms on both sides gives

$$\log p(a_\varepsilon | \text{rest}) = -2 \log(a_\varepsilon) - (\sigma_\varepsilon^{-2} + A_\varepsilon^{-2}) / a_\varepsilon + \text{const},$$

which then leads to the full conditional

$$a_\varepsilon | \text{rest} \sim \text{Inverse-Gamma}(1, \sigma_\varepsilon^{-2} + A_\varepsilon^{-2}).$$

Taking expectations with respect to all parameters except  $a_\varepsilon$  gives

$$\begin{aligned}\log q^*(a_\varepsilon) &= -2 \log(a_\varepsilon) - E_q(\sigma_\varepsilon^{-2} + A_\varepsilon^{-2}) / a_\varepsilon + \text{const} \\ &= (-1 - 1) \log(a_\varepsilon) - (\mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}) / a_\varepsilon + \text{const}.\end{aligned}$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

The expressions for  $B_{q(a_\varepsilon)}$  and  $\mu_{q(1/a_\varepsilon)}$  follow immediately.

**Expressions for  $q^*(\sigma_w^2)$ ,  $B_{q(\sigma_w^2)}$  and  $\mu_{q(1/\sigma_w^2)}$**

$$q^*(\sigma_w^2) \sim \text{Inverse-Gamma}\left(\frac{1}{2}(\sum_{i=1}^m n_i + 1), B_{q(\sigma_w^2)}\right),$$

where

$$B_{q(\sigma_w^2)} = \frac{1}{2} E_q(\|\mathbf{w} - \mathbf{x}\|^2) + \mu_{q(1/\sigma_w^2)} \quad \text{and} \quad \mu_{q(\sigma_w^2)} = \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_w^2)}$$

and

$$\begin{aligned} E_q(\|\mathbf{w} - \mathbf{x}\|^2) &= \|\mathbf{w}_{x_{\text{obs}}} - \mathbf{x}_{\text{obs}}\|^2 + \|\mathbf{w}_{x_{\text{unobs}}}\|^2 - 2\mathbf{w}_{x_{\text{unobs}}}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})} \\ &\quad + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}\|^2 + \text{tr}(\text{Cov}_q(\mathbf{x}_{\text{unobs}})) \\ &= \|\mathbf{w}_{x_{\text{unobs}}} - \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}\|^2 + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2. \end{aligned}$$

*Derivation:*

The derivation is analogous to that of  $q^*(\sigma_\varepsilon^2)$  and related quantities.

**Expressions for  $q^*(a_w)$ ,  $B_{q(a_w)}$  and  $\mu_{q(1/a_w)}$**

$$q^*(a_w) \sim \text{Inverse-Gamma}(1, B_{q(a_w)}),$$

where

$$B_{q(a_w)} = \mu_{q(1/\sigma_w^2)} + A_w^{-2} \quad \text{and} \quad \mu_{q(1/a_w)} = 1 / \{\mu_{q(1/\sigma_w^2)} + A_w^{-2}\}.$$

*Derivation:*

The derivation is analogous to that of  $q^*(a_\varepsilon)$  and related quantities.

**Expressions for  $q^*(\mu_x)$ ,  $B_{q(\mu_x)}$  and  $\mu_{q(1/\mu_x)}$**

$$q^*(\mu_x) \sim N(\mu_{q(\mu_x)}, \sigma_{q(\mu_x)}^2)$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

where

$$\mu_{q(\mu_x)} = \sigma_{q(\mu_x)}^2 \mu_{q(1/\sigma_x^2)} (\mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}) \quad \text{and} \quad \sigma_{q(\mu_x)}^2 = 1/\{N\mu_{q(1/\sigma_x^2)} + 1/\sigma_{\mu_x}^2\}.$$

*Derivation:*

First we note that

$$p(\mu_x | \text{rest}) \propto p(\mu_x | \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_x^2) \propto p(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}} | \mu_x, \sigma_x^2) p(\mu_x).$$

Taking logarithms on both sides gives

$$\begin{aligned} \log p(\mu_x | \text{rest}) &= -\frac{1}{2\sigma_x^2} \|\mathbf{x} - \mu_x \mathbf{1}\|^2 - \frac{1}{2\sigma_{\mu_x}^2} \mu_x^2 \\ &= -\frac{1}{2} \left\{ \left( \frac{N}{\sigma_x^2} + \frac{1}{\sigma_{\mu_x}^2} \right) \mu_x^2 - 2\frac{1}{\sigma_x^2} \mu_x \mathbf{1}^\top \mathbf{x} \right\}. \end{aligned}$$

Taking expectations with respect to all parameters except  $\mu_x$  gives

$$\log p(\mu_x | \text{rest}) = -\frac{1}{2} \left\{ (N\mu_{q(1/\sigma_x^2)} + 1/\sigma_{\mu_x}^2) \mu_x^2 - 2\mu_{q(1/\sigma_x^2)} \mu_x E_q(\mathbf{1}^\top \mathbf{x}) \right\}.$$

The form of  $q^*(\sigma_x^2)$  and  $B_{q(\sigma_x^2)}$  follows where

$$E_q(\mathbf{1}^\top \mathbf{x}) = \mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}.$$

**Expressions for  $q^*(\sigma_x^2)$ ,  $B_{q(\sigma_x^2)}$  and  $\mu_{q(1/\sigma_x^2)}$**

$$q^*(\sigma_x^2) \sim \text{Inverse-Gamma} \left( \frac{1}{2} (\sum_{i=1}^m n_i + 1), B_{q(\sigma_x^2)} \right),$$

where

$$B_{q(\sigma_x^2)} = \frac{1}{2} E_q(\|\mathbf{x} - \mu_x \mathbf{1}\|^2) + \mu_{q(1/\sigma_x^2)} \quad \text{and} \quad \mu_{q(\sigma_x^2)} = \frac{1}{2} (\sum_{i=1}^m n_i + 1) / B_{q(\sigma_x^2)}$$

and

$$\begin{aligned} E_q(\|\mathbf{x} - \mu_x \mathbf{1}\|^2) &= \|\mathbf{x}_{\text{obs}} - \mu_{q(\mu_x)} \mathbf{1}\|^2 + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})} - \mu_{q(\mu_x)} \mathbf{1}\|^2 \\ &\quad + N\sigma_{q(\mu_x)}^2 + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2. \end{aligned}$$

*Derivation:*

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

The derivation is analogous to that of  $q^*(\sigma_\varepsilon^2)$  and related quantities.

The remaining  $q^*$  densities involve straightforward adaptation of the derivations given in Appendix 2.A.

##### 4.A.1 Derivation of the marginal log-likelihood lower bound

The expression for the marginal log-likelihood lower bound is

$$\begin{aligned}
\log \underline{p}(\mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{w}; q) &= -\frac{1}{2}\sigma_\beta^{-2} \left( \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right) + \frac{1}{2} \left( \sum_{\ell=1}^L q_\ell^G + p + m \right) \\
&+ \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| - \log(\mathcal{C}_{q^R, \nu+q^R-1}) + \log(\mathcal{C}_{q^R, \nu+m+q^R-1}) \\
&- \frac{1}{2} (\nu + m + q^R - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma}^R)}| - \frac{1}{2} \sum_{i=1}^m n_i^{\text{unobs}} \log(\sigma_{q(\mathbf{x}_{\text{unobs}, i})}^2) \\
&- \frac{1}{2} \sum_{\ell=1}^L (q_\ell^G + 1) \log(B_{q(\sigma_{u_\ell}^2)}) - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_\varepsilon^2)}) \\
&- \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_x^2)}) - \frac{1}{2} (\sum_{i=1}^m n_i + 1) \log(B_{q(\sigma_w^2)}) \\
&- \sum_{\ell=1}^L \log(B_{q(a_{u_\ell})}) + \sum_{\ell=1}^L \mu_{q(1/a_{u_\ell})} \mu_{q(1/\sigma_{u_\ell}^2)} \\
&+ \sum_{r=1}^{q^R} \nu(\mathbf{M}_{q((\boldsymbol{\Sigma}^R)^{-1})_{rr}}) \mu_{q(1/a_r^R)} - \frac{1}{2} (\nu + q^R) \sum_{r=1}^{q^R} \log(B_{q(a_r^R)}) \\
&- \log(B_{q(a_\varepsilon)}) + \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)} - \log(B_{q(a_x)}) + \mu_{q(1/a_x)} \mu_{q(1/\sigma_x^2)} \\
&- \log(B_{q(a_w)}) + \mu_{q(1/a_w)} \mu_{q(1/\sigma_w^2)} \\
&- \frac{1}{2} \left\{ \log(\sigma_{q(\mu_x)}^2) - \log(\sigma_{\mu_x}^2) \right\} - \frac{1}{2} \sigma_{\mu_x}^{-2} \|\mu_{q(\mu_x)}\|^2 \\
&+ \frac{1}{2} \sigma_{q(\mu_x)}^2 \left( 1/\sigma_{q(\mu_x)}^2 - 1/\sigma_{\mu_x}^2 \right) + \text{const},
\end{aligned}$$

where  $\ell = 1$ .

*Derivation:*

The logarithm of the lower bound on the marginal likelihood is given by

$$\begin{aligned}
\log \underline{p}(\mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{w}; q) &= \\
&E_q \left\{ \log p(\mathbf{y}, \mathbf{x}, \mathbf{w}, \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \boldsymbol{\Sigma}^R, \sigma_u^2, \mathbf{a}^R, \mathbf{a}_u, \mu_x, \sigma_x^2, a_x, \sigma_w^2, a_w, \sigma_\varepsilon^2, a_\varepsilon) \right. \\
&\quad \left. - \log q^*(\mathbf{x}, \mathbf{w}, \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \boldsymbol{\Sigma}^R, \sigma_u^2, \mathbf{a}^R, \mathbf{a}_u, \mu_x, \sigma_x^2, a_x, \sigma_w^2, a_w, \sigma_\varepsilon^2, a_\varepsilon) \right\}.
\end{aligned}$$

According to the product restriction (4.5), our expression for the lower bound of the log-likelihood becomes:

$$\begin{aligned}
\log \underline{p}(\mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{w}; q) &= E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_\varepsilon^2) \right\} \\
&+ E_q \left\{ \log p(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G | \boldsymbol{\Sigma}^R, \sigma_u^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G) \right\} \\
&+ E_q \left\{ \log p(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}} | \mu_x, \sigma_x^2) - \log q^*(\mathbf{x}_{\text{unobs}}) \right\}
\end{aligned}$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

$$\begin{aligned}
& + E_q \{ \log p(\boldsymbol{\Sigma}^R | \mathbf{a}^R) - \log q^*(\boldsymbol{\Sigma}^R) \} E_q \{ p(\sigma_u^2 | \mathbf{a}_u) - \log q^*(\sigma_u^2) \} \\
& + E_q \{ p(\sigma_x^2 | a_x) - \log q^*(\sigma_x^2) \} + E_q \{ p(\sigma_w^2 | a_w) - \log q^*(\sigma_w^2) \} \\
& + E_q \{ \log p(\mathbf{a}^R) - \log q^*(\mathbf{a}^R) \} + E_q \{ \log p(\mathbf{a}_u) - \log q^*(\mathbf{a}_u) \} \\
& + E_q \{ \log p(a_x) - \log q^*(a_x) \} + E_q \{ \log p(a_w) - \log q^*(a_w) \} \\
& + E_q \{ p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \} + E_q \{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \} \\
& + E_q \{ \log p(\mathbf{w}_{x_{\text{obs}}}, \mathbf{w}_{x_{\text{unobs}}} | \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_w^2) \} \\
& + E_q \{ \log p(\mu_x) - \log q^*(\mu_x) \}.
\end{aligned}$$

We are only interested in the ones highlighted in darker colour as they represent the differences between models (2.11) and (4.1).

First note that

$$\log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_\varepsilon^2) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2.$$

Taking expectations gives

$$\begin{aligned}
E_q \{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}^R, \mathbf{u}^G, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}}, \sigma_\varepsilon^2) \} = \\
-\frac{N}{2} \log(2\pi) - \frac{N}{2} E_q \{ \log(\sigma_\varepsilon^2) \} - \frac{1}{2\sigma_\varepsilon^2} E_q \{ \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 \},
\end{aligned}$$

where

$$\begin{aligned}
E_q \{ \|\mathbf{y} - \mathbf{C}\mathbf{v}\|^2 \} = & \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}_{x_{\text{obs}}}^\top \mathbf{C}_{x_{\text{obs}}} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \\
& + \|\mathbf{y}_{x_{\text{unobs}}}\|^2 - 2 \mathbf{y}_{x_{\text{unobs}}}^\top E_q(\mathbf{C}_{x_{\text{unobs}}}) \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \\
& + \text{tr} \left\{ E_q(\mathbf{C}_{x_{\text{unobs}}}^\top \mathbf{C}_{x_{\text{unobs}}}) (\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^\top + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \right\}.
\end{aligned}$$

The next term is

$$\begin{aligned}
\log p(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}} | \mu_x, \sigma_x^2) - \log q^*(\mathbf{x}_{\text{unobs}}) = \\
-\frac{1}{2} \left\{ \sum_{i=1}^m (n_i \log(2\pi) - n_i \log(\sigma_x^2) - \sigma_x^{-2} \|\mathbf{x}_i - \mu_x \mathbf{1}\|^2) \right\} \\
-\frac{1}{2} \left\{ \sum_{i=1}^m (n_i^{\text{unobs}} \log(2\pi) - n_i^{\text{unobs}} \log(\sigma_{q(\mathbf{x}_{\text{unobs}}, i)}^2)) \right. \\
\left. - (1/\sigma_{q(\mathbf{x}_{\text{unobs}}, i)}^2) \|\mathbf{x}_{\text{unobs}, i} - \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}}, i)}\|^2 \right\}.
\end{aligned}$$

#### 4.A. OPTIMAL $q$ -DENSITIES DERIVATION FOR MEASUREMENT ERROR PROBLEMS

---

Taking expectations gives

$$\begin{aligned}
E_q \{ \log p(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{unobs}} | \mu_x, \sigma_x^2) - \log q^*(\mathbf{x}_{\text{unobs}}) \} = \\
\sum_{i=1}^m \left\{ -\frac{1}{2} n_i^{\text{obs}} \log(2\pi) - \frac{1}{2} n_i^{\text{obs}} E_q(\log(\sigma_x^2)) - \frac{1}{2} n_i^{\text{unobs}} \log(\sigma_{q(\mathbf{x}_{\text{unobs},i})}^2) \right. \\
\left. - \frac{1}{2} \mu_{q(1/\sigma_x^2)} \left( \|\mathbf{x}_{\text{obs},i} - \mu_{q(\mu_x)} \mathbf{1}_{n_i^{\text{obs}}}\|^2 + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs},i})} - \mu_{q(\mu_x)} \mathbf{1}_{n_i^{\text{unobs}}}\|^2 + n_i \sigma_{q(\mu_x)}^2 \right) \right. \\
\left. + \frac{1}{2} \left( 1/\sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 + \mu_{q(1/\sigma_x^2)} \right) n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2 \right\}.
\end{aligned}$$

The next term is

$$\begin{aligned}
\log p(\sigma_x^2 | a_x) - \log q^*(\sigma_x^2) &= -\frac{1}{2} \log(a_x) - \log \Gamma(1/2) - \frac{3}{2} \log(\sigma_x^2) - 1/(a_x \sigma_x^2) \\
&\quad - \frac{1}{2} (N+1) \log(B_{q(\sigma_x^2)}) + \log \Gamma\{\frac{1}{2}(N+1)\} \\
&\quad + \left(\frac{N}{2} + \frac{3}{2}\right) \log(\sigma_x^2) - B_{q(\sigma_x^2)}/\sigma_x^2.
\end{aligned}$$

Taking expectations gives

$$\begin{aligned}
E_q \{ \log p(\sigma_x^2 | a_x) - \log q^*(\sigma_x^2) \} &= -\frac{1}{2} E_q \{ \log(a_x) \} - \log \Gamma(1/2) + \frac{N}{2} E_q \{ \log(\sigma_x^2) \} \\
&\quad + \mu_{q(1/\sigma_x^2)} (B_{q(\sigma_x^2)} - \mu_{q(1/a_x)}) \\
&\quad - \frac{1}{2} (N+1) \log B_{q(\sigma_x^2)} + \log \Gamma\{\frac{1}{2}(N+1)\}.
\end{aligned}$$

Similar arguments lead to

$$\begin{aligned}
E_q \{ \log p(\sigma_w^2 | a_w) - \log q^*(\sigma_w^2) \} &= -\frac{1}{2} E_q \{ \log(a_w) \} - \log \Gamma(1/2) + \frac{N}{2} E_q \{ \log(\sigma_w^2) \} \\
&\quad + \mu_{q(1/\sigma_w^2)} (B_{q(\sigma_w^2)} - \mu_{q(1/a_w)}) \\
&\quad - \frac{1}{2} (N+1) \log B_{q(\sigma_w^2)} + \log \Gamma\{\frac{1}{2}(N+1)\}
\end{aligned}$$

and

$$\begin{aligned}
E_q \{ \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \} &= -\frac{1}{2} E_q \{ \log(a_\varepsilon) \} - \log \Gamma(1/2) + \frac{N}{2} E_q \{ \log(\sigma_\varepsilon^2) \} \\
&\quad + \mu_{q(1/\sigma_\varepsilon^2)} (B_{q(\sigma_\varepsilon^2)} - \mu_{q(1/a_\varepsilon)}) \\
&\quad - \frac{1}{2} (N+1) \log B_{q(\sigma_\varepsilon^2)} + \log \Gamma\{\frac{1}{2}(N+1)\}.
\end{aligned}$$

The next term is

$$E_q \{ \log p(\mathbf{w} | \mathbf{x}, \sigma_w^2) \} = -\frac{N}{2} \log(2\pi) - \frac{N}{2} E_q \{ \log(\sigma_w^2) \} - \frac{1}{2} \mu_{q(1/\sigma_w^2)} E_q(\|\mathbf{w} - \mathbf{x}\|^2),$$

where

$$\begin{aligned} E_q(\|\mathbf{w} - \mathbf{x}\|^2) &= \|\mathbf{w}_{x_{\text{obs}}} - \mathbf{x}_{\text{obs}}\|^2 + \|\mathbf{w}_{x_{\text{unobs}}}\|^2 - 2\mathbf{w}_{x_{\text{unobs}}}^\top \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})} \\ &\quad + \|\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}\|^2 + \text{tr}(\text{Cov}_q(\mathbf{x}_{\text{unobs}})) \\ &= \|\mathbf{w}_{x_{\text{unobs}}} - \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}\|^2 + \sum_{i=1}^m n_i^{\text{unobs}} \sigma_{q(\mathbf{x}_{\text{unobs},i})}^2. \end{aligned}$$

The next term is

$$\begin{aligned} \log p(\mu_x) - \log q^*(\mu_x) &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_{\mu_x}^2) - \frac{1}{2\sigma_{\mu_x}^2} \|\mu_x\|^2 \\ &\quad + \frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\sigma_{q(\mu_x)}^2) + \frac{1}{2\sigma_{q(\mu_x)}^2} (\|\mu_x - \mu_{q(\mu_x)}\|^2). \end{aligned}$$

Taking expectations gives

$$\begin{aligned} E_1 \{\log p(\mu_x) - \log q^*(\mu_x)\} &= -\frac{1}{2} \left( \log(\sigma_{q(\mu_x)}^2) - \log(\sigma_{\mu_x}^2) \right) - \frac{1}{2\sigma_{\mu_x}^2} \|\mu_{q(\mu_x)}\|^2 \\ &\quad + \frac{1}{2} \sigma_{q(\mu_x)}^2 \left( \frac{1}{\sigma_{q(\mu_x)}^2} - \frac{1}{\sigma_{\mu_x}^2} \right). \end{aligned}$$

The remaining expectation expressions follow similar algebraic calculations and are therefore omitted.

## 4.B Optimal $q$ -densities derivation for missing data problems

We first write out the full model expression

$$y_{ij} = \beta_0 + \beta_x x_{ij} + \beta_s s_{ij} + \sum_{k=1}^{q^G} u_k^G z_k(s_{ij}) + u_{0i}^R + u_{1i}^R x_{ij} + \varepsilon_{ij},$$

which can be broken into the observed and missing parts:

$$\begin{bmatrix} \mathbf{y}_{x_{\text{obs}}} \\ \mathbf{y}_{x_{\text{mis}}} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{n_{\text{obs}}} & \mathbf{x}_{\text{obs}} & \mathbf{s}_{x_{\text{obs}}} \\ \mathbf{1}_{n_{\text{mis}}} & \mathbf{x}_{\text{mis}} & \mathbf{s}_{x_{\text{mis}}} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_x \\ \beta_s \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_{x_{\text{obs}}}^G \\ \mathbf{Z}_{x_{\text{mis}}}^G \end{bmatrix} \begin{bmatrix} u_1^G \\ \vdots \\ u_m^G \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_{x_{\text{obs}}}^R \\ \mathbf{Z}_{x_{\text{mis}}}^R \end{bmatrix} \begin{bmatrix} \mathbf{u}_1^R \\ \vdots \\ \mathbf{u}_m^R \end{bmatrix}.$$

**Expression for  $q^*(\mathbf{x}_{\text{mis},i})$ ,  $\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})}$  and  $\sigma_{q(\mathbf{x}_{\text{mis},i})}^2$**

For  $i = 1, \dots, m$

$$q^*(\mathbf{x}_{\text{mis},i}) \sim N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})}, \sigma_{q(\mathbf{x}_{\text{mis},i})}^2 \mathbf{I}),$$

where

$$\begin{aligned} \sigma_{q(\mathbf{x}_{\text{mis},i})}^2 &= E_q \left\{ \mathbf{1}_{n_i^{\text{mis}}}^\top \left( \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^{\text{R}})^2 \mathbf{1}_{n_i^{\text{mis}}} + \frac{1}{\sigma_x^2} \mathbf{1}_{n_i^{\text{mis}}} - 2\phi_1^2 \lambda(\boldsymbol{\xi}_{x_{\text{mis},i}}) \right) \right\}^{-1} \\ \text{and } \boldsymbol{\mu}_{q(\mathbf{x}_{\text{mis},i})} &= \sigma_{q(\mathbf{x}_{\text{mis},i})}^2 \mathbf{I}_{x_{\text{mis}}} E_q \left\{ \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^{\text{R}})^2 \left( \mathbf{y}_{x_{\text{mis},i}} - \beta_0 \mathbf{1}_{n_i^{\text{mis}}} \right. \right. \\ &\quad \left. \left. - \beta_s \mathbf{s}_{x_{\text{mis},i}} - \sum_{k=1}^{q^{\text{G}}} u_k^{\text{G}} z_k(\mathbf{s}_{x_{\text{mis},i}}) - u_{0i}^{\text{R}} \mathbf{1}_{n_i^{\text{mis}}} \right) \right. \\ &\quad \left. + \left( \frac{\mu_x}{\sigma_x^2} + 2\phi_0 \phi_1 \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis},i}})\} - \frac{1}{2}\phi_1 \right) \mathbf{1}_{n_i^{\text{mis}}} \right\} \end{aligned}$$

*Derivation:*

First we note that

$$\begin{aligned} p(\mathbf{x}_{\text{mis}}|\text{rest}) &\propto p(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}}, \mathbf{y}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \boldsymbol{\beta}, \sigma_\varepsilon^2, \mu_x, \sigma_x^2, \mathbf{R}, \boldsymbol{\phi}) \\ &\propto p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}, \mathbf{u}^{\text{R}}, \mathbf{u}^{\text{G}}, \sigma_\varepsilon^2) p(\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}|\mu_x, \sigma_x^2, \mathbf{R}, \boldsymbol{\phi}) p(\mathbf{R}|\mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}, \boldsymbol{\phi}). \end{aligned}$$

Taking logarithms on both sides where  $x$  is missing gives

$$\begin{aligned} \log p(\mathbf{x}_{\text{mis}}|\text{rest}) &= \sum_{i=1}^m \left[ -\frac{1}{2\sigma_\varepsilon^2} \left\| \mathbf{y}_{x_{\text{mis},i}} - \left( \beta_0 \mathbf{1}_{n_i^{\text{mis}}} + \beta_x \mathbf{x}_{\text{mis},i} + \beta_s \mathbf{s}_{x_{\text{mis},i}} \right. \right. \right. \\ &\quad \left. \left. - \sum_{k=1}^{q^{\text{G}}} u_k^{\text{G}} z_k(\mathbf{s}_{x_{\text{mis},i}}) + u_{0i}^{\text{R}} \mathbf{1}_{n_i^{\text{mis}}} + u_{1i}^{\text{R}} \mathbf{x}_{\text{mis},i} \right\|^2 \right. \\ &\quad \left. - \frac{1}{2\sigma_x^2} \left\| \mathbf{x}_{\text{mis},i} - \mu_x \mathbf{1}_{n_i^{\text{mis}}} \right\|^2 + \{1 + \exp(\phi_0 + \phi_1 \mathbf{x}_{\text{mis},i})\}^{-1} \right] \\ &\geq -\frac{1}{2} \left( \mathbf{1}_{n_i^{\text{mis}}}^\top \left[ \left\{ \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^{\text{R}})^2 \mathbf{1}_{n_i^{\text{mis}}} + \frac{1}{\sigma_x^2} \mathbf{1}_{n_i^{\text{mis}}} - 2\phi_1^2 \lambda\{\boldsymbol{\xi}_{x_{\text{mis},i}}\} \right\} \odot (\mathbf{x}_{\text{mis},i})^2 \right] \right. \\ &\quad \left. - 2 \mathbf{x}_{\text{mis},i}^\top \left\{ \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^{\text{R}}) \left( \mathbf{y}_{x_{\text{mis},i}} - \beta_0 \mathbf{1}_{n_i^{\text{mis}}} - \beta_s \mathbf{1}_{n_i^{\text{mis}}} - \sum_{k=1}^{q^{\text{G}}} u_k^{\text{G}} z_k(\mathbf{s}_{x_{\text{mis},i}}) - u_{0i}^{\text{R}} \mathbf{1}_{n_i^{\text{mis}}} \right) \right. \right. \right. \\ &\quad \left. \left. + \frac{1}{\sigma_x^2} (\mu_x \mathbf{1}_{n_i^{\text{mis}}}) + \left( 2\phi_0 \phi_1 \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis},i}})\} - \frac{1}{2}\phi_1 \right) \mathbf{1}_{n_i^{\text{mis}}} \right\} \right) + \text{const}, \end{aligned}$$

where the two expectation terms are as follows:

$$E_q \left\{ \mathbf{1}_{n_i^{\text{mis}}}^\top \left( \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^{\text{R}})^2 \mathbf{1}_{n_i^{\text{mis}}} + \frac{1}{\sigma_x^2} \mathbf{1}_{n_i^{\text{mis}}} - 2\phi_1^2 \lambda(\boldsymbol{\xi}_{x_{\text{mis},i}}) \right) \right\}$$



$$\begin{aligned}
 &= \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \boldsymbol{\mu}_{q(\beta_x)}^2 + \boldsymbol{\Sigma}_{q(\beta_x)} + 2(\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{1i}^R)}) \right. \\
 &\quad \left. + \boldsymbol{\mu}_{q(u_{1i}^R)}^2 + \boldsymbol{\Sigma}_{q(u_{1i}^R)} \right\} + \mu_{q(1/\sigma_x^2)} \\
 &\quad + 2(\boldsymbol{\mu}_{q(\phi_1)}^2 + \boldsymbol{\Sigma}_{q(\phi_1)}) \sum_{i=1}^m \lambda(\boldsymbol{\xi}_{x_{\text{mis},i}}), \\
 E_q \left\{ \frac{1}{\sigma_\varepsilon^2} (\beta_x + u_{1i}^R)^2 \left( \mathbf{y}_{x_{\text{mis},i}} - \beta_0 \mathbf{1}_{n_i^{\text{mis}}} - \beta_s \mathbf{s}_{x_{\text{mis},i}} - \sum_{k=1}^{q_G} u_k^G z_k(\mathbf{s}_{x_{\text{mis},i}}) - u_{0i}^R \mathbf{1}_{n_i^{\text{mis}}} \right) \right. \\
 &\quad \left. + \left( \frac{\mu_x}{\sigma_x^2} + 2\phi_0 \phi_1 \text{diag}\{\lambda(\boldsymbol{\xi}_{x_{\text{mis},i}})\} - \frac{1}{2}\phi_1 \right) \mathbf{1}_{n_i^{\text{mis}}} \right\} \\
 &= \sigma_{q(x_{\text{mis},i})} \mathbf{I}_{n_i^{\text{mis}}} \left[ \mu_{q(1/\sigma_\varepsilon^2)} \left\{ (\boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\mu}_{q(u_{1i}^R)}) \mathbf{y}_{x_{\text{mis},i}} \right. \right. \\
 &\quad - (\boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(\beta_x)} + \boldsymbol{\Lambda}_{q(\beta_0, \beta_x)} + \boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_0, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{mis}}} \\
 &\quad - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(\beta_s)} + \boldsymbol{\Lambda}_{q(\beta_x, \beta_s)} + \boldsymbol{\mu}_{q(\beta_s)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(\beta_s, u_{1i}^R)}) \mathbf{s}_{x_{\text{mis},i}} \\
 &\quad - \sum_{k=1}^{q_G} (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\beta_x, u_k^G)}) z_k(\mathbf{s}_{x_{\text{mis},i}}) \\
 &\quad - \sum_{k=1}^{q_G} (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(\beta_x, u_k^G)}) z_k(\mathbf{s}_{x_{\text{mis},i}}) \\
 &\quad - \sum_{k=1}^{q_G} (\boldsymbol{\mu}_{q(u_{1i}^R)} \boldsymbol{\mu}_{q(u_k^G)} + \boldsymbol{\Lambda}_{q(u_{1i}^R, u_k^G)}) z_k(\mathbf{s}_{x_{\text{mis},i}}) \\
 &\quad \left. \left. - (\boldsymbol{\mu}_{q(\beta_x)} \boldsymbol{\mu}_{q(u_{0i}^R)} + \boldsymbol{\Lambda}_{q(\beta_x, u_{0i}^R)} + \boldsymbol{\mu}_{q(u_{0i}^R)} \boldsymbol{\mu}_{q(u_{1i}^R)} + \boldsymbol{\Lambda}_{q(u_{0i}^R, u_{1i}^R)}) \mathbf{1}_{n_i^{\text{mis}}} \right\} \right. \\
 &\quad \left. + \mu_{q(1/\sigma_x^2)} \boldsymbol{\mu}_{q(\mu_x)} \mathbf{1}_{n_i^{\text{mis}}} + 2(\boldsymbol{\mu}_{q(\phi_1)}^2 + \boldsymbol{\Sigma}_{q(\phi_1)}) \lambda(\boldsymbol{\xi}_{x_{\text{mis},i}}) \right].
 \end{aligned}$$

Expressions for  $q^*(\phi)$ ,  $\boldsymbol{\mu}_{q(\phi)}$  and  $\boldsymbol{\Sigma}_{q(\phi)}$

$$q^*(\phi) \sim N(\boldsymbol{\mu}_{q(\phi; \boldsymbol{\xi})}, \boldsymbol{\Sigma}_{q(\phi; \boldsymbol{\xi})})$$

where

$$\boldsymbol{\mu}_{q(\phi; \boldsymbol{\xi})} = \boldsymbol{\Sigma}_{q(\phi; \boldsymbol{\xi})} \left\{ (\mathbf{R} - \frac{1}{2}) E_q(\mathbf{X}^\top) \right\} \quad \text{and} \quad \boldsymbol{\Sigma}_{q(\phi; \boldsymbol{\xi})} = \left( \frac{1}{\sigma_\phi^2} \mathbf{I} - 2 E_q[\mathbf{X}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{X}^\top] \right)^{-1}.$$

*Derivation:*

First we note that

$$\begin{aligned}
 p(\phi | \text{rest}) &\propto p(\phi | \mathbf{R}, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}) = p(\mathbf{R} | \phi, \mathbf{x}_{\text{obs}}, \mathbf{x}_{\text{mis}}) p(\phi) \\
 &\propto \exp[\mathbf{R}^\top (\mathbf{X} \phi) - \mathbf{1}^\top \log(1 + \exp(\mathbf{X} \phi))] \exp\left(-\frac{1}{2\sigma_\phi^2} \|\phi\|^2\right)
 \end{aligned}$$

Replacing

$$-\mathbf{1}^\top \log(1 + \exp(\mathbf{X} \phi)) \geq (\mathbf{X} \phi)^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} (\mathbf{X} \phi) - \frac{1}{2} \mathbf{1}^\top (\mathbf{X} \phi) + \mathbf{1}^\top \psi(\boldsymbol{\xi})$$

and taking logarithms on both sides gives

$$\log p(\boldsymbol{\phi}|\text{rest}) \geq -\frac{1}{2} \left\{ \boldsymbol{\phi}^\top \left( \frac{1}{\sigma_\phi^2} \mathbf{I} - 2 [\mathbf{X}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{X}] \right) \boldsymbol{\phi} - 2 \boldsymbol{\phi}^\top (\mathbf{R} - \frac{1}{2}) \mathbf{X}^\top \right\}.$$

Taking expectations with respect to all parameters except  $\mathbf{x}_{\text{unobs},i}$  leads to

$$\log q^*(\boldsymbol{\phi}) \geq -\frac{1}{2} \left\{ \boldsymbol{\phi}^\top \left( \frac{1}{\sigma_\phi^2} \mathbf{I} - 2 E_q[\mathbf{X}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{X}] \right) \boldsymbol{\phi} - 2 \boldsymbol{\phi}^\top (\mathbf{R} - \frac{1}{2}) E_q(\mathbf{X})^\top \right\},$$

where the expressions for  $E_q[\mathbf{X}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{X}]$  and  $E_q(\mathbf{X})$  are given in Subsection 4.3.1.

## Chapter 5

# Extension to Group-Specific Curve Models With Contrasting

*There is a magic in graphs. The profile of a curve reveals in a flash a whole situation - the life history of an epidemic, a panic, or an era of prosperity. The curve informs the mind, awakens the imagination, convinces.*

Henry D. Hubbard

### 5.1 Introduction

Parametric regression methods for longitudinal and multilevel data have been well developed in the last two decades. Such methods can be broadly classified into estimating equations based methods (Diggle *et al.*, 2002; Zeger and Liang, 1986) and mixed models (Breslow and Clayton, 1993; Gelman and Hill, 2007; Goldstein, 2011; Laird and Ware, 1982). A parametric regression model assumes that the relationship between the mean of a response outcome and predictor variables to be of a known functional form. Although such a parametric assumption offers simplicity, it is inappropriate for situations when the relationship between the mean response and predictors is unknown. As an illustration, Figure 5.1 shows the overall and hospital-specific trends in caesarean section rates for low-risk nulliparous women in the largest state of Australia, New South Wales (NSW),

---

The full content of this chapter is published as: Lee, C. Y. Y. and Wand, M. P. (2015). Variational methods for fitting complex Bayesian mixed models to health data. *Statistics in Medicine*. DOI: 10.1002/sim.6737. A manuscript related to this chapter is currently under revision in *Birth*: Lee, C. Y. Y., Homer, C., Bisits, A. and Ryan, L. (2015). Increasing variation in hospital caesarean section rates among low-risk nulliparous women in Australia, from 1994 to 2010.

between 1994 and 2010. The trends are clearly non-linear and show substantial within- and between-hospital variability over time. In order to capture these features simultaneously, a more complicated and flexible model is required. For example, a more flexible model is a generalisation to additive models (Hastie and Tibshirani, 1990), where the overall mean and deviation of each group (in this case hospital) from that overall mean are modelled nonparametrically with arbitrary smooth functions of time. We refer to such model as a *group-specific curve model*.

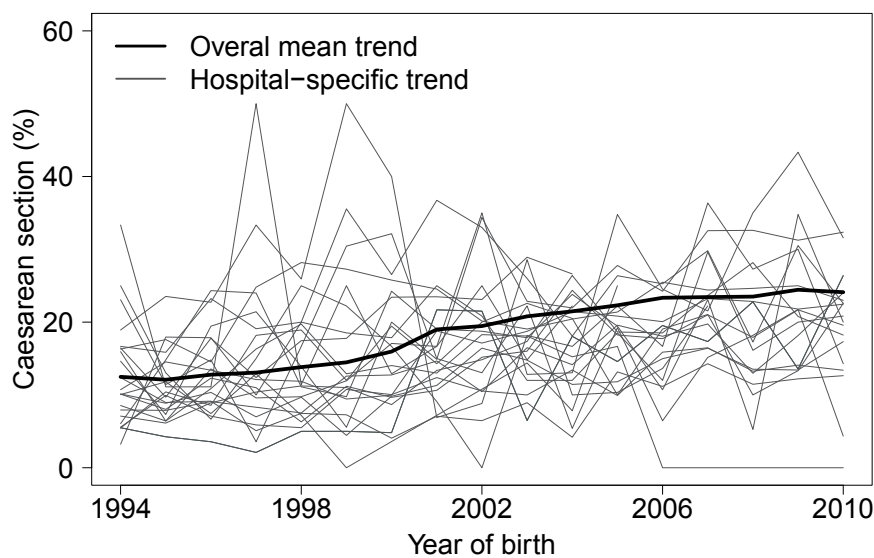


Figure 5.1: Trends in caesarean section rates for low-risk nulliparous women in the largest state of Australia, New South Wales, between 1994 and 2010. The dotted line is the overall mean trend and the solid lines are the selected hospital-specific trends.

Early work on group-specific curve models used different smoothing techniques (e.g. kernels and smoothing splines) to estimate both overall mean curve and group-specific curves, with random effects modelled either by a parametric function or a Gaussian process (Verbyla *et al.*, 1999; Zeger and Diggle, 1994; Zhang *et al.*, 1998). Brumback and Rice (1998) embedded smoothing splines within the mixed model framework. However, they ran into computational problems in the estimation since they assumed fixed intercepts and slopes for the group-specific curves, thus leading to the number of fixed effects being at least twice as large as the number of curves. Rice and Wu (2001) alleviated computational problems by modelling group-specific curves as spline basis functions with random coefficients. However, their approach requires a careful selection of the number and location of knots. Durbán *et al.* (2005) relaxed the importance of this selection by using low rank smoothing splines with a penalty approach, which expresses the group-specific curves as a

linear combination of truncated polynomial spline bases with random coefficients. Their proposed covariance matrix for the random basis coefficients was modelled parametrically with independence assumptions. Chen and Wang (2011) and Ryu *et al.* (2011) extended this work by allowing a more general covariance matrix structure for the random basis coefficients, although the former was confined to functional data analysis.

When the dimension of the spline basis functions and the number of groups become large, many of the above-mentioned approaches become too slow, or even computationally infeasible. These problems motivate the development of fast and scalable methods for fitting group-specific curve models to large longitudinal and multilevel datasets. In Bayesian statistics, exact inference for nonparametric or semiparametric regression models that use penalised spline basis functions is typically intractable, requiring approximate inference methods for use in practice. Markov chain Monte Carlo or MCMC is the most commonly used approximate inference method in this setting, but can be computationally intensive and often suffers from poor convergence in large and complex models. A faster, deterministic alternative to MCMC is *variational approximations* (Bishop, 2006; Ormerod and Wand, 2010). The basic idea behind variational approximations is to recast the problem of computing posterior probabilities as an optimisation problem by introducing a class of more manageable approximating distributions, and then optimizing some criterion to find the distribution within that class that best matches the posterior. Recently, there has been a growing interest in variational methods for longitudinal and multilevel models. Ormerod and Wand (2010) and Luts *et al.* (2014) developed variational algorithms for fitting and inference for grouped data. However, their algorithms are infeasible for datasets with a large number of groups due to naïve inversion of the random effects covariance matrix. Tan and Nott (2013) introduced a partially noncentered nonparametrisation strategy allowing the random effects for each group to be independent. Such a restriction was shown to improve efficiency of the variational algorithms. Lee and Wand (2015a) did not impose such a restriction, but rather, took advantage of the inherent blocked structure and sparseness of the random effects covariance matrix and developed algorithms that streamline its inversion and estimation. Their streamlined algorithms for the longitudinal and multilevel models result in impressive speed improvements (up to the order of thousands) and currently represent the state-of-the-art in this area. Stewart (2014) applied their streamlined algorithms to research studies in social sciences.

In this chapter, we present a fully variational approach to fitting a series of Bayesian logistic mixed models in order to characterise trends in caesarean section rates in New South Wales, Australia. We begin with the standard random intercept and slope model and then generalize the model by replacing linear models for the overall mean and hospital-specific trends with arbitrary smooth functions that are nonparametrically estimated from the data. Section 5.2 reviews the random intercept and slope model and presents various

extensions of the standard model including nonparametric functions and factor-by-curve interactions. A basic framework of approximate inference for a simple Bayesian logistic mixed model is outlined, serving as a building block for the more complicated models. Details on variational fitting and inference for models are also given. In Section 5.3 we present numerical evidence of the efficacy of our developed variational algorithms in terms of inferential accuracy and computational speed. The final proposed model is illustrated in Section 5.4 with the analysis of caesarean section data. The chapter concludes with a brief discussion in Section 5.5.

## 5.2 Mean field variational Bayes approximations to caesarean section data

In this section, we present an epidemiology study that motivates our methodological development. Caesarean section rates are increasing worldwide and among countries of Organisation for Economic Co-operation and Development the average rate increased from 20% in 2000 to 27% in 2011 (Organisation for Economic Co-Operation and Development, 2011). It is recognised that caesarean section rates vary considerably across regions and hospitals in several countries. For example, in the United States, a four-fold variation was found between low- and high-use areas (Baicker *et al.*, 2006). In the United Kingdom, rates of emergency cesarean section among National Health Service trusts ranged from 15% to 32% (Bragg *et al.*, 2010). In the largest state of Australia, NSW, there were 1 500 964 deliveries from 1994 to 2010, with caesarean section rate increased from 17.4% to 30.6% over this 17-year period. While the statewide trend in caesarean section rate is well reported, little is known about the trend in caesarean section rate for each hospital, especially for small hospitals. Previous work by the author addressed this gap by examining hospital-specific trends in caesarean section rates for low-risk nulliparous women in NSW between 1994 and 2010 (Lee *et al.*, 2015). Here we extend the model applied previously, using the MFVB approximation, to compare hospital-specific trends between subpopulations of women.

Data were obtained from the NSW Perinatal Data Collection. The Perinatal Data Collection is a legislated population-based surveillance system covering all live births, and stillbirths of at least 20 weeks gestation or at least 400 grams birthweight. Information is recorded by the attending midwife or medical practitioner providing maternity care, and includes maternal demographic, medical and obstetric information of the mother, and details of labor, birth and condition of the infant. Data for this analysis were de-identified. Permission for use of data was approved by the NSW Ministry of Health. The study population consists of low-risk women giving birth for the first time in a NSW public

or private hospital between 1994 and 2010. Low-risk women were defined as those women giving birth for the first time (nulliparous women), aged 20-34 years who did not smoke in pregnancy and did not have any pre-existing or gestational medical conditions, such as diabetes and hypertension, gave birth to singleton cephalic (head down position) live infants at term (37 weeks gestation).

The primary outcome is the annual rate of caesarean section for each hospital in NSW over a 17-year period. Public and private hospitals with  $\geq 50$  births per annum for more than half of the 17-year study period were included. The data have a two-level hierarchical structure with women nested within hospitals. To define notation, we use a double subscript convention with  $i$  denoting hospital ( $i = 1, \dots, m$ ) and  $j$  denoting woman ( $j = 1, \dots, n_i$ ). For example,  $y_{ij}$  denote the binary indicator of caesarean section for woman  $j$  in hospital  $i$  and  $x_{ij}$  the corresponding year of birth. The total number of low-risk nulliparous women in the study population is  $\sum_{i=1}^m n_i$ . To examine the underlying trends in caesarean section rates overall and for each hospital, we consider a series of Bayesian logistic mixed models with increasing complexity. For each model, the year of birth is modelled by adding together two **arbitrary unspecified functions** as covariates, one representing the overall mean  $f(x)$  and the other being hospital-specific departures from that overall mean  $g_i(x)$ . To keep notation simple, we consider only a single predictor in the model. Extension to models with more than one predictor is straightforward. More precisely, we assume

$$y_{ij} \stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\text{logit}^{-1}\{f(x_{ij}) + g_i(x_{ij})\}), \quad 1 \leq j \leq n_i, \quad 1 \leq i \leq m, \quad (5.1)$$

where the functional forms of  $f$  and  $g_i$  will be reviewed in the later subsections. The function  $g_i$  controls for hospital-to-hospital variation, both in magnitude and across the 17-year trend. The sum of  $f$  and  $g_i$  provides the trend in caesarean section rate for each hospital, which can have an unique intercept, slope and even shape depending upon whether they are of a parametric or nonparametric form. The term  $f(x_{ij}) + g_i(x_{ij})$  represents the log odds of caesarean section for woman  $j$  in hospital  $i$  and will be the focus in the following subsections.

The Bayesian hierarchical model corresponding to (5.1) is

$$\mathbf{y}|\boldsymbol{\beta}, \mathbf{u} \sim \text{Bernoulli}(\text{logit}^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})), \quad \mathbf{u}|\mathbf{G} \sim N(\mathbf{0}, \mathbf{G}). \quad (5.2)$$

The notation  $\mathbf{v} \sim \text{Bernoulli}(\mathbf{p})$  used in (5.2) is shorthand for the entries of the random vector  $\mathbf{v}$  having independent Bernoulli distributions with parameters corresponding to the entries of  $\mathbf{p}$ . The matrices  $\mathbf{X}$  and  $\mathbf{Z}$  are the respective  $\sum_{i=1}^m n_i \times P$  fixed effects design matrix and  $\sum_{i=1}^m n_i \times d$  random effects design matrix,  $\boldsymbol{\beta}$  and  $\mathbf{u}$  are the so called  $P \times 1$

fixed effects vector and  $d \times 1$  random effects vector. The random effects vectors are given a multivariate normal prior with zero mean and covariance matrix  $\mathbf{G}$ . Throughout we take the prior distribution of the fixed effects vector  $\boldsymbol{\beta}$  to be of the form  $\boldsymbol{\beta} \sim N(0, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I})$ , where  $\sigma_{\boldsymbol{\beta}}^2$  is to be chosen by the analyst. In logistic regression, situations where a shift in the predictor  $x$  corresponds to the probability of response  $y$  changing from 0.01 to 0.99 are rare. Hence, a prior distribution that assigns low probabilities to changes of 10 on the logistic scale would be appropriate. Further, we use proper but “diffuse” conditionally conjugate priors for the random effects, this corresponds to a Half-Cauchy distribution for a single variance component (Result 5 of Wand *et al.* (2012))

$$\sigma^2 | a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/a), \quad a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/A^2).$$

The multivariate extension of the Half-Cauchy distribution is a scaled inverse-Wishart distribution for an unstructured  $q^R \times q^R$  random effects covariance matrix (Huang and Wand, 2013)

$$\begin{aligned} \boldsymbol{\Sigma}_R | a_{R,1}, \dots, a_{R,q^R} &\sim \text{Inverse-Wishart}(\nu + q^R - 1, 2\nu \text{diag}(1/a_{R,1}, \dots, 1/a_{R,q^R})), \\ a_{R,r} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}(\frac{1}{2}, 1/A_R^2), \quad 1 \leq r \leq q^R, \end{aligned} \tag{5.3}$$

with hyperparameters  $\nu, A, A_R > 0$ . The value  $\nu = 2$  corresponds to the correlation parameters having uniform distributions over  $(-1, 1)$  and the standard deviation parameters having Half- $t$  distributions with 2 degrees of freedom.

It will soon become apparent that the model framework introduced in (5.2) encompasses a wide range of logistic mixed models with different types of group structures by simply changing  $\mathbf{G}$ .

### 5.2.1 Random intercept and slope model

A natural starting point for these caesarean section data is to assume that the overall mean and deviation of the  $i$ th hospital from that overall mean are simply straight lines. This leads to the standard random intercept and slope model (Laird and Ware, 1982), with  $f(x)$  and  $g_i(x)$  modelling through **linear functions**:

$$\begin{aligned} f(x) &= \beta_0 + \beta_1 x_{ij} \quad ; \quad g_i(x) = U_{0i} + U_{1i} x_{ij} \\ \text{and } \mathbf{G} &= \mathbf{I}_m \otimes \boldsymbol{\Sigma}_R = \text{blockdiag}_{1 \leq i \leq m} \begin{bmatrix} \sigma_{u_0}^2 & \sigma_{u_0, u_1} \\ \sigma_{u_0, u_1} & \sigma_{u_1}^2 \end{bmatrix}, \end{aligned} \tag{5.4}$$

where  $\beta_0$  and  $\beta_1$  are the respective overall intercept and slope, and  $U_{0i}$  and  $U_{1i}$  are the hospital-specific deviations from that overall intercept and slope, being treated as a random



## 5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

---

sample from a bivariate normal distribution with an unstructured  $2 \times 2$  covariance matrix  $\Sigma_{\mathbf{R}}$ . The model accounts for possible variability in the intercepts  $\sigma_{u_0}^2$  and slopes  $\sigma_{u_1}^2$  of hospitals and allows for an intercept-slope correlation  $\rho_{u_0, u_1}$ . In common practice, we generalise the fixed and random components of the mixed models to arbitrary general design matrices (Section 2 of Zhao *et al.* (2006)). This allows one to take advantage of the ever-expanding methods and software for inference in these models. Henceforth, we rewrite model (5.4) in matrix notation as

$$\begin{aligned} \mathbf{X}_i &\equiv \begin{bmatrix} 1 & x_{i1} \\ \vdots & \vdots \\ 1 & x_{in_i} \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{bmatrix}, \quad \boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \\ \mathbf{Z} &\equiv \begin{bmatrix} \mathbf{X}_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{X}_m \end{bmatrix} \quad \text{and} \quad \mathbf{u} \equiv \begin{bmatrix} u_{01} \\ u_{11} \\ \vdots \\ u_{0m} \\ u_{1m} \end{bmatrix}. \end{aligned} \tag{5.5}$$

Noting that the random effects design matrix  $\mathbf{Z}$  has a block-diagonal form, where each block corresponds to the  $i$ th hospital-specific deviations from the overall mean and slope.

We now combine the ideas introduced in Chapter 1 Subsection 1.5.3 to develop a scalable iterative algorithm for approximate Bayesian inference in model (5.4). We first briefly introduce the following useful notation: for a scalar random variable  $\theta$ , let  $\mu_{q(\theta)} \equiv E_q(\theta)$  and  $\sigma_{q(\theta)} \equiv \text{Var}_q(\theta)$  be the mean and variance with respect to the approximating  $q$ -distribution. For a random vector parameter  $\boldsymbol{\theta}$ , we use the analogously defined  $\boldsymbol{\mu}_{q(\boldsymbol{\theta})}$  and  $\Sigma_{q(\boldsymbol{\theta})}$ . In addition, we define  $\mathbf{M}_{q(\boldsymbol{\Theta})}$  to be the mean with respect to the approximating  $q$ -distribution for a random matrix  $\boldsymbol{\Theta}$ . The journey towards a practical MFVB algorithm commences with a  $q$ -density product restriction

$$p(\boldsymbol{\beta}, \mathbf{u}, \mathbf{G}, a_{\mathbf{R},1}, a_{\mathbf{R},2} | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\mathbf{G}) q(a_{\mathbf{R},1}) q(a_{\mathbf{R},2}),$$

where  $a_{\mathbf{R},1}$  and  $a_{\mathbf{R},2}$  are auxiliary parameters defined in (5.3). Such a restriction arises from an interaction between the initial factorisation assumed in the approximating posterior and the underlying conditional independence properties of the true joint posterior (Bishop, 2006).

To provide an example of how optimal  $q$ -densities are constructed, we derive the optimal density  $q^*(\boldsymbol{\beta}, \mathbf{u})$ . According to Algorithm 2, the optimal  $q$ -density for  $(\boldsymbol{\beta}, \mathbf{u})$  satisfies

$$q^*(\boldsymbol{\beta}, \mathbf{u}) \propto E_q\{\log p(\boldsymbol{\beta}, \mathbf{u}, \mathbf{y}, \mathbf{G}, a_{\mathbf{R},1}, a_{\mathbf{R},2})\} \tag{5.6}$$

$$\begin{aligned} &\propto E_q\{\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}) \times \log p(\mathbf{u}|\mathbf{G}) \times \log p(\boldsymbol{\beta})\} \\ &\propto E_q[\mathbf{y}^\top(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^\top \log\{\mathbf{1} + \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})\} - \frac{1}{2}\mathbf{u}^\top \mathbf{G}^{-1}\mathbf{u} - \frac{1}{2\sigma_\beta^2}\boldsymbol{\beta}^\top \boldsymbol{\beta}]. \end{aligned}$$

We see that the non-quadratic convex term  $-\log(1 + e^x)$  in the likelihood poses a multivariate intractability problem with regard to approximate inference for  $(\boldsymbol{\beta}, \mathbf{u})$ . To get around this, we transform this convex term to a simple quadratic function  $f(x; \xi)$ , a trick first introduced in Jaakkola and Jordan (1997). Different values of  $\xi$  correspond to different parabolas, all of which are smaller than  $-\log(1 + e^x)$ . Thus we can simplify the convex term to be the maxima of a family of parabolas (Jaakkola and Jordan, 1997)

$$-\log(1 + e^x) = \max_{\xi \in \mathbb{R}} \left\{ \lambda(\xi) x^2 - \frac{1}{2}x + \psi(\xi) \right\} \quad \text{for all } x \in \mathbb{R}, \quad (5.7)$$

where  $\lambda(\xi) \equiv -\tanh(\xi/2)/(4\xi)$  and  $\psi(\xi) \equiv \xi/2 - \log(1 + e^\xi) + \xi \tanh(\xi/2)/4$ . Let  $\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$  and  $\mathbf{v} \equiv [\boldsymbol{\beta}^\top \ \mathbf{u}^\top]^\top$  and substitute (5.7) into (5.6) gives the following lower bound on  $q^*(\boldsymbol{\beta}, \mathbf{u})$ :

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}; \xi) &\geq \mathbf{y}^\top \mathbf{C}\mathbf{v} - \left[ \mathbf{v}^\top \mathbf{C}^\top \text{diag}\{\lambda(\xi)\} \mathbf{C}\mathbf{v} - \frac{1}{2}\mathbf{1}^\top \mathbf{C}\mathbf{v} \right] - \frac{1}{2}\mathbf{v}^\top \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & E_q(\mathbf{G}^{-1}) \end{bmatrix} \mathbf{v} \\ &= -\frac{1}{2} \left\{ \mathbf{v}^\top \left( 2\mathbf{C}^\top \text{diag}\{\lambda(\xi)\} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & E_q(\mathbf{G}^{-1}) \end{bmatrix} \right) \mathbf{v} - \mathbf{C}^\top \left( \mathbf{y} - \frac{1}{2}\mathbf{1} \right) \right\} \\ &= \log \underline{q}^*(\boldsymbol{\beta}, \mathbf{u}; \xi), \end{aligned} \quad (5.8)$$

where  $E_{q(\mathbf{G}^{-1})} \equiv \mathbf{I}_m \otimes \mathbf{M}_{q(\boldsymbol{\Sigma}_R^{-1})}$  and  $\boldsymbol{\xi}$  is being introduced as an  $\sum_{i=1}^m n_i \times 1$  vector of variational parameters. Scalar functions applied to vectors are evaluated element-wise. Noting that since the right hand side of (5.8) is a quadratic form and, by completing the square in the usual way to identify the mean and covariance, we approximate the posterior of  $(\boldsymbol{\beta}, \mathbf{u})$  by a multivariate normal distribution

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) &\sim N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}), \\ \text{where } \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} &\equiv \left( 2\mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & E_q(\mathbf{G}^{-1}) \end{bmatrix} \right)^{-1} \\ \text{and } \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} &\equiv \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \mathbf{C}^\top \left( \mathbf{y} - \frac{1}{2}\mathbf{1} \right). \end{aligned}$$

The remaining  $q$ -densities involve similar adaptation of the above derivation and hence we

## 5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

---

omit the details and state their functional forms directly:

$$\begin{aligned}\boldsymbol{\xi} &\leftarrow \sqrt{\text{diagonal}\{\mathbf{C}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^\top) \mathbf{C}^\top\}}, \\ q^*(\boldsymbol{\Sigma}_{\mathbf{R}}) &\text{ is the Inverse-Wishart } (\nu + m + 1, \mathbf{B}_{q(\boldsymbol{\Sigma}_{\mathbf{R}})}) \text{ density function,} \\ q^*(a_{\mathbf{R}, r}) &\text{ is the Inverse-Gamma } (1, B_{q(a_{\mathbf{R}, r})}) \text{ density function, } 1 \leq r \leq 2,\end{aligned}$$

where the parameter  $\mathbf{B}_{q(\boldsymbol{\Sigma}_{\mathbf{R}})}$  is the scale matrix of the inverse-Wishart  $q$ -density and  $B_{q(a_{\mathbf{R}, r})}$  is the scale parameter of the inverse-Gamma  $q$ -density. Thus all optimal densities belong to parametric families with the parameters explicitly determined by the distributions of the remaining model parameters and observed data. Taken together, these solutions lead to Algorithm 15 for approximate Bayesian inference in model (5.4).

---

**Initialise:**  $\boldsymbol{\xi}$  ( $\sum_{i=1}^m n_i \times 1$ ; all entries positive),  $\mathbf{M}_{q(\boldsymbol{\Sigma}_{\mathbf{R}}^{-1})}$  positive definite and  $\mu_{q(a_{\mathbf{R}, r})} > 0$ ,  $r = 1, 2$

**Cycle through updates:**

**Define:**  $E_{q(\mathbf{G}^{-1})} \equiv \mathbf{I}_m \otimes \mathbf{M}_{q(\boldsymbol{\Sigma}_{\mathbf{R}}^{-1})}$

**Update mean and covariance matrix of multivariate normal  $q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$ :**

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \leftarrow \left( 2 \mathbf{C}^\top \text{diag}\{\lambda(\boldsymbol{\xi})\} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & E_{q(\mathbf{G}^{-1})} \end{bmatrix} \right)^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \mathbf{C}^\top (\mathbf{y} - \frac{1}{2} \mathbf{1})$$

$$\boldsymbol{\xi} \leftarrow \sqrt{\text{diagonal}\{\mathbf{C}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^\top) \mathbf{C}^\top\}}$$

**Update mean and scale matrix of inverse-Wishart  $q^*(\boldsymbol{\Sigma}_{\mathbf{R}})$ :**

$$\mathbf{B}_{q(\boldsymbol{\Sigma}_{\mathbf{R}})} \leftarrow \sum_{i=1}^m (\mu_{q(\mathbf{u}_i)} \mu_{q(\mathbf{u}_i)} + \boldsymbol{\Sigma}_{q(\mathbf{u}_i)}) + 2\nu \text{diag}(\mu_{q(1/a_{\mathbf{R}, 1})}, \mu_{q(a_{\mathbf{R}, 2})})$$

$$\mathbf{M}_{q(\boldsymbol{\Sigma}_{\mathbf{R}}^{-1})} \leftarrow \frac{1}{2} (\nu + m + 1) \mathbf{B}_{q(\boldsymbol{\Sigma}_{\mathbf{R}})}^{-1}$$

**Update mean and scale parameter of inverse-Gamma  $q^*(a_{\mathbf{R}, r})$ :**

For  $r = 1, 2$ :

$$B_{q(a_{\mathbf{R}, r})} \leftarrow \nu (\mathbf{M}_{q(\boldsymbol{\Sigma}_{\mathbf{R}}^{-1})})_{rr} + 1/A_{\mathbf{R}}^2 \quad ; \quad \mu_{q(a_{\mathbf{R}, r})} \leftarrow \frac{1}{2} (\nu + 1) / B_{q(a_{\mathbf{R}, r})}$$

**until the increase in  $p(\mathbf{y}; q)$  is negligible.**

---

Algorithm 15: Iterative scheme for obtaining the optimal  $q$ -density functions in the Bayesian logistic mixed model (5.4).

Algorithm 15 is a basic framework for approximate inference for a simple Bayesian logistic mixed model that serves as a building block for the more complicated models. In what follows, we will gradually increase the complexity of model (5.4) and demonstrate a

great advantage of the MFVB algorithms, *modularity*. By this we mean that concepts like main effects, interaction effects, higher-order random effects and spline regression can be viewed as modules that can be put together into an almost endless variety of statistical models. All that required is relatively straightforward modifications on the structures of the general design matrices and random effects covariance matrix in order to accommodate larger and more complicated mixed models.

### 5.2.2 Group-specific curve model

Close inspection of Figure 5.1 shows that the straight line assumption imposed by the random intercept and slope model is unreasonable. We seek to relax this linearity assumption by replacing the linear mean and hospital-specific functions with **arbitrary smooth functions**, hereafter group-specific curve model, taking form (5.1), but with

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^{K_{\text{gbl}}} u_{\text{gbl},k} z_{\text{gbl},k}(x) ; \quad g_i(x) = U_{0i} + U_{1i} x + \sum_{k=1}^{K_{\text{grp}}} u_{\text{grp},ik} z_{\text{grp},ik}(x)$$

$$\text{and } \mathbf{G} = \begin{bmatrix} \overbrace{\sigma_{\text{gbl}}^2 \mathbf{I}_{K_{\text{gbl}}}}^{\text{Overall mean}} & \overbrace{\mathbf{0}}^{\text{Hospital-specific}} \\ \mathbf{0} & \text{blockdiag}_{1 \leq i \leq m} (\Sigma_{\text{R}} | \sigma_{\text{grp}}^2 \mathbf{I}_{K_{\text{grp}}}) \end{bmatrix}. \quad (5.9)$$

This model is a direct extension of the random intercept and slope model since in (5.9) the overall mean curve and each hospital-specific curve have two components: a linear part (analogous to model (5.4)) and a non-linear part, which allows more flexibility.

There are numerous approaches to modelling and estimating  $f$  and  $g_i$  nonparametrically. The one which is most conducive to inference via variational methods is penalised regression splines with mixed model representation. An advantage of this approach is that one can think in terms of constructing a series of basis functions that can be used as covariates, such as  $z_{\text{gbl},k}$  and  $z_{\text{grp},ik}$  being the respective spline bases of size  $K_{\text{gbl}}$  and  $K_{\text{grp}}$ . Our preference of  $z_{\text{gbl},k}$  and  $z_{\text{grp},ik}$  is suitably linearly transformed cubic O'Sullivan penalised splines (Section 4 of Wand and Ormerod (2008)), since this leads to approximate smoothing splines with good boundary and extrapolation properties. In practice,  $K_{\text{gbl}} = 25$  is a sufficient choice for most spline basis functions (Li and Ruppert, 2008) and typically  $K_{\text{grp}}$  is smaller than  $K_{\text{gbl}}$  since fewer basis functions are needed to handle group-specific deviations. The coefficients  $u_{\text{gbl},k}$  and  $u_{\text{grp},ik}$  can be considered as a measure of the basis amplitude since they regulate the roughness of the time curves. In order to avoid overfitting the data, we impose a penalty on the basis coefficients by treating them as a random sample from a normal distribution with mean 0 and variance  $\sigma^2$ , i.e.  $u_{\text{gbl},k} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{\text{gbl}}^2)$  and  $u_{\text{grp},ik} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{\text{grp}}^2)$  respectively. The variances  $\sigma_{\text{gbl}}^2$  and  $\sigma_{\text{grp}}^2$  are

## 5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

---

commonly known as the *smoothing parameters*, thus selection of smoothing parameters in model (5.9) simply reduces to variance component estimation in a mixed effect model.

From an extendability standpoint, a general design framework for mixed effect models is a particularly attractive approach that allows smoothing-type models such as group-specific curve models to be fitted as generalized linear mixed models. Moving from parametric regression to nonparametric regression using mixed model based penalised splines involves replacing  $\mathbf{Z}$

$$\text{from } \begin{array}{c} \text{Hospital 1} \quad \text{Hospital } m \\ \begin{bmatrix} \mathbf{X}_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{X}_m \end{bmatrix} \end{array} \quad \text{to} \quad \begin{array}{c} \text{Overall mean} \quad \text{Hospital 1} \quad \cdots \quad \text{Hospital } m \\ \begin{bmatrix} \mathbf{X}_1 | \mathbf{Z}_{\text{grp},1} & \cdots & \mathbf{0} \\ \mathbf{Z}_{\text{gbl}} & \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{X}_m | \mathbf{Z}_{\text{grp},m} \end{bmatrix} \end{array}, \quad (5.10)$$

where  $\mathbf{Z}_{\text{gbl}}$  is an  $\sum_{i=1}^m n_i \times K_{\text{gbl}}$  random spline basis matrix for the overall mean curve and  $\mathbf{Z}_{\text{grp},i}$  is an  $n_i \times K_{\text{grp}}$  random spline basis matrix for the  $i$ th hospital-specific deviation. Further, even with the additional  $K_{\text{gbl}} + mK_{\text{grp}}$  random effects parameters, the covariance matrix  $\mathbf{G}$  shown in (5.9) remains to be a block-diagonal, but certainly larger, matrix with different entries for the variance components corresponding to the random basis coefficients for the overall mean and random effects for each hospital. Taking these into consideration, it then follows that the MFVB algorithm for model (5.9) simply involves modifications of Algorithm 16 to incorporate the above-mentioned structural changes in the  $\mathbf{Z}$  and  $\mathbf{G}$  matrices, as well as updates for the additional model parameters as follows:

$$\begin{aligned} q^*(a_{\text{gbl}}) & \text{ is the Inverse-Gamma } \left(1, B_{q(a_{\text{gbl}})}\right) \text{ density function;} \\ q^*(\sigma_{\text{gbl}}^2) & \text{ is the Inverse-Gamma } \left(\frac{1}{2}(K_{\text{gbl}} + 1), B_{q(\sigma_{\text{gbl}}^2)}\right) \text{ density function;} \\ q^*(a_{\text{grp}}) & \text{ is the Inverse-Gamma } \left(1, B_{q(a_{\text{grp}})}\right) \text{ density function;} \text{ and} \\ q^*(\sigma_{\text{grp}}^2) & \text{ is the Inverse-Gamma } \left(\frac{1}{2}(mK_{\text{grp}} + 1), B_{q(\sigma_{\text{grp}}^2)}\right) \text{ density function,} \end{aligned}$$

where  $B_{q(a_{\text{gbl}})}, B_{q(a_{\text{grp}})}, B_{q(\sigma_{\text{gbl}}^2)}$  and  $B_{q(\sigma_{\text{grp}}^2)}$  are the scale parameters with respect to the corresponding  $q$ -density functions.

### 5.2.3 Factor-by-curve interactions

The second objective of this caesarean section analysis is to examine differences in trends in caesarean section rates between low-risk nulliparous women aged less than 25 years and those aged greater than or equal to 25 years. We are therefore interested in fitting a separate time curve for each maternal age group (see Figure 5.2). This leads to a group-specific curve model with a *factor-by-curve interaction*, where one can view the effect of

## 5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

maternal age on the probability of caesarean section varies smoothly over time. From here onwards we refer low-risk nulliparous women aged less than 25 years as Group *A* (younger group) and those greater or equal to 25 years as Group *B* (older group).

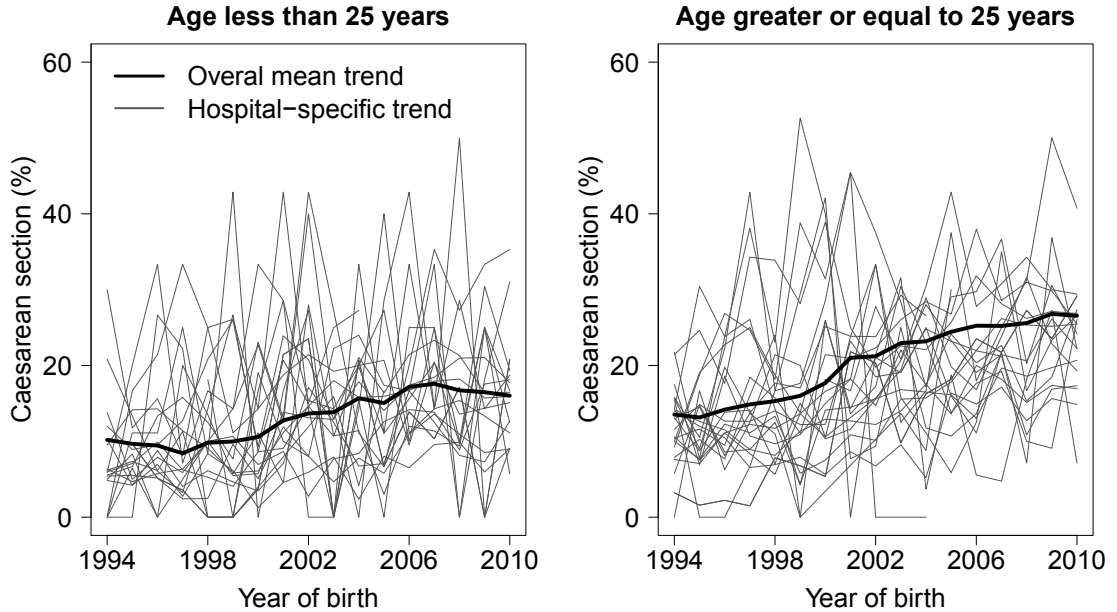


Figure 5.2: Trends in caesarean section rates for low-risk nulliparous women aged less than 25 years and those aged greater or equal to 25 years in New South Wales, Australia, from 1994 to 2010. The dotted line is the overall mean trend and the solid lines are the selected hospital-specific trends.

Additive models that include factor-by-curve interactions are known to be useful but can be computationally very demanding. Coull *et al.* (2000) used a backfitting algorithm, in which, for each backfit iteration corresponding to a curve interaction term, one splits the data into subsets and fits a smooth function to each subset. However, this backfitting approach can be computationally impractical when the number of interaction terms or the number of data subsets become large. Maringwa *et al.* (2008) used penalised regression splines in a mixed model framework to compare population-averaged profiles, but only focused on the random intercept models without any group-specific curves. The penalised spline approach is preferable since the former data subsetting approach must be nested within a backfitting algorithm. We follow on Maringwa *et al.* (2008) and incorporate a factor-by-curve interaction into the group-specific curve model (5.9). Let  $I_{ij}^A = 1$  if  $(x_{ij}, y_{ij})$  belongs to a low-risk nulliparous women of aged less than 25 years and zero otherwise, we propose a model of the following form

$$(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_{ij} = I_{ij}^A \{f^A(x_{ij}) + g_i^A(x_{ij})\} + (1 - I_{ij}^A) \{f^B(x_{ij}) + g_i^B(x_{ij})\},$$

## 5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

---

where

$$f^A(x) = \beta_0^A + \beta_1^A x + \sum_{k=1}^{K_{\text{gbl}}} u_{\text{gbl},k}^A z_{\text{gbl},k}(x) \quad (5.11)$$

$$g_i^A(x) = U_{0i}^A + U_{1i}^A x \sum_{k=1}^{K_{\text{grp}}} u_{\text{grp},ik}^A z_{\text{grp},k}(x),$$

$$f^B(x) = \beta_0^A + \beta_0^{\text{BvsA}} + (\beta_1^A + \beta_1^{\text{BvsA}}) x + \sum_{k=1}^{K_{\text{gbl}}} u_{\text{gbl},k}^B z_{\text{gbl},k}(x)$$

$$g_i^B(x) = U_{0i}^B + U_{1i}^B x \sum_{k=1}^{K_{\text{grp}}} u_{\text{grp},ik}^B z_{\text{grp},k}(x)$$

$$\text{and } \mathbf{G} = \begin{bmatrix} \overbrace{(\sigma_{\text{gbl}}^A)^2 \mathbf{I}_{K_{\text{gbl}}}}^{\text{Overall mean}} & \mathbf{0} & \overbrace{\mathbf{0}}^{\text{Hospital-specific}} \\ \mathbf{0} & (\sigma_{\text{gbl}}^B)^2 \mathbf{I}_{K_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag} \left( \sum_{\text{R}} | \sigma_{\text{grp}}^2 \mathbf{I}_{K_{\text{grp}}} | \sigma_{\text{grp}}^2 \mathbf{I}_{K_{\text{grp}}} \right)_{1 \leq i \leq m} \end{bmatrix}.$$

We consider low-risk women aged less than 25 years as our reference group and define the corresponding  $f^A$  and  $g_i^A$  functions the same way as in model (5.9). For low-risk women aged greater or equal to 25 years, the  $f^B$  function is defined slightly differently by introducing the parameters  $\beta_0^{\text{BvsA}}$  and  $\beta_1^{\text{BvsA}}$  that represent the differences in overall intercepts and slopes between low-risk nulliparous women of younger and older age. This gives an *asymmetrical structure* in the fixed effects formulation. While asymmetrical formulation of fixed effects is common in the presence of a reference group, we do not recommend such a structure imposed on the random spline coefficients since it would induce restriction on the smoothness of a particular curve. We impose a different variance parameter for each overall mean curve, allowing the level of smoothing to differ between curves. A simpler model would be to assume a common variance parameter for all overall mean curves, i.e.  $\sigma_{\text{gbl}}^A = \sigma_{\text{gbl}}^B = \sigma_{\text{gbl}}$ , meaning that all curves have equivalent smoothness, but with different forms of shape. In addition, we account for the potential correlation between the random intercepts and slopes of each of the women age groups within the same hospital using an unstructured  $4 \times 4$  covariance matrix, i.e.  $[U_{0i}^A \ U_{1i}^A \ U_{0i}^B \ U_{1i}^B]^\top \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \mathbf{\Sigma}_{\text{R}})$ .

The differences in trends in caesarean section rates for low-risk nulliparous women aged greater or equal to 25 years versus those aged less than 25 years can be quantified via estimation of a so-called *contrast function*. At the population level, this is simply obtained

by subtracting the two overall mean curves:

$$c^{\text{BvsA}}(x) \equiv f^{\text{B}}(x) - f^{\text{A}}(x) = \beta_0^{\text{BvsA}} + \beta_1^{\text{BvsA}} x + \sum_{k=1}^{K_{\text{gbl}}} (u_{\text{gbl},k}^{\text{B}} - u_{\text{gbl},k}^{\text{A}}) z_{\text{gbl},k}(x). \quad (5.12)$$

This contrast function can be interpreted as the log odds of caesarean section (averaged across hospitals) for low-risk nulliparous women aged greater or equal to 25 years, as a function of time, compared with those aged less than 25 years.

The structural changes in the random effects matrices  $\mathbf{Z}$  and  $\mathbf{G}$  underpin the complexity of a mixed effect model being transitioned from the simplest random intercept and slope model to a more complicated model including a factor-by-curve interaction. For example, Figure 5.3 shows that the  $\mathbf{Z}$  matrix starts off as having a simple block-diagonal structure and, as the model increases in complexity, it grows into a “nested” block-diagonal structure (Each large block is itself block-diagonal, with one or more small blocks on the diagonal). As the number of random spline basis functions increases in the model, the dimensions of the  $\mathbf{Z}$  and  $\mathbf{G}$  matrices increase accordingly. While these changes do not directly affect the algebraic aspect of Algorithm 16 because of the advantage of modularity, one must be aware that the update of the covariance matrix  $\Sigma_{q(\beta, \mathbf{u}; \xi)}$  now requires storage and inversion of a large and sparse matrix. Specifically, the number of columns in  $\Sigma_{q(\beta, \mathbf{u}; \xi)}$  for each considered model is

1. Random intercepts and slopes  $P + m q^{\text{R}}$
2. Group-specific curves  $P + K_{\text{gbl}} + m (q^{\text{R}} + K_{\text{grp}})$
3. Factor-by-curve interactions  $P + 2K_{\text{gbl}} + 2m (q^{\text{R}} + K_{\text{grp}})$

Without doubt the number of groups  $m$  dominates the dimension of  $\Sigma_{q(\beta, \mathbf{u}; \xi)}$ . For example, if  $P = 10, m = 1000, K_{\text{gbl}} = 25, q^{\text{R}} = 2$  and  $K_{\text{grp}} = 10$ , then the dimension of the matrix requiring storage and inversion for the most complicated model would be  $24\,060 \times 24\,060$ . This aspect renders Algorithm 16 infeasible for very large and complex longitudinal and multilevel datasets. As pointed out by Lee and Wand (2015a), the naïve implementation of MFVB algorithms for arbitrarily large grouped data is extremely inefficient in terms of speed and storage. The time complexity can be as high as  $O(m^3)$ , making variational inference impractical for large and complex mixed models. Through exploiting the inherent block-diagonal structure of the random effects covariance matrix, the authors developed fast and memory-efficient MFVB algorithms that streamline inversion and update for  $\Sigma_{q(\beta, \mathbf{u}; \xi)}$ . Here we briefly reiterate their streamlined approach in more general terms. Recall that the general update expression for  $\Sigma_{q(\beta, \mathbf{u}; \xi)}$  in Algorithm



5.2. MEAN FIELD VARIATIONAL BAYES APPROXIMATIONS TO CAESAREAN SECTION DATA

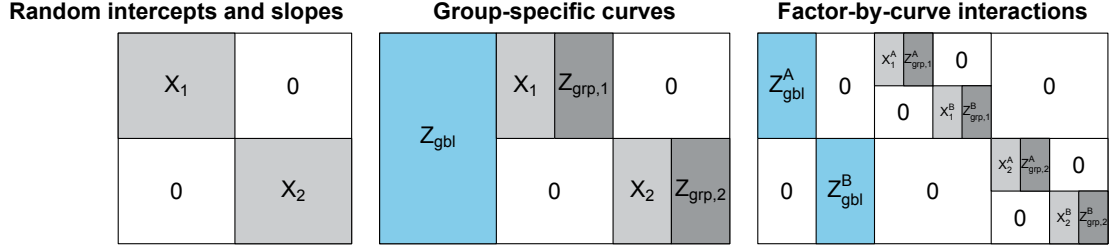


Figure 5.3: Various structures of the  $\mathbf{Z}$  matrix across logistic mixed models (5.4), (5.9) and (5.11), with the number of groups  $m = 2$ . The  $\mathbf{Z}$  matrix starts off as having a simple block-diagonal structure and, as the model increases in complexity, it grows into a “nested” block-diagonal structure. The definitions of matrices are described in Subsection 5.2.2.

15 involves inversion of the following matrix

$$2 \mathbf{C}^\top \text{diag} \{ \lambda(\boldsymbol{\xi}) \} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & E_{q(\mathbf{G}^{-1})} \end{bmatrix} = \left[ \begin{array}{c|cccc} \mathbf{U} & \mathbf{V}_1 & \mathbf{V}_2 & \cdots & \mathbf{V}_m \\ \hline \mathbf{V}_1^\top & \mathbf{W}_1^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{V}_2^\top & \mathbf{0} & \mathbf{W}_2^{-1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{V}_m^\top & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{W}_m^{-1} \end{array} \right] \quad (5.13)$$

The central idea of the streamlined approach requires permuting the above matrix into an approximate block-diagonal form for decomposition, as shown in (5.13). The matrices  $\mathbf{U}$ ,  $\mathbf{V}_i$  and  $\mathbf{W}_i^{-1}$  are usually of general forms with small dimension (Appendix of (Lee and Wand, 2015a)); each can be easily derived via straightforward matrix manipulation. With this transformation we can efficiently invert the block-partitioned form of (5.13) by solving a system of simultaneous linear equations in matrix form (Smith and Wand, 2008)

$$\begin{bmatrix} \overbrace{\mathbf{U} \quad \mathbf{V}}^{\Sigma_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}^{-1}} \\ \overbrace{\mathbf{V}^\top \quad \mathbf{W}} \end{bmatrix} \begin{bmatrix} \overbrace{\mathbf{X} \quad \mathbf{Y}}^{\Sigma_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}} \\ \overbrace{\mathbf{Y}^\top \quad \mathbf{Z}} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

Standard calculations lead to the following forms of submatrices:

$$\begin{aligned} \mathbf{X} &\equiv (\mathbf{U} - \mathbf{V} \mathbf{W}^{-1} \mathbf{V}^\top)^{-1} \\ \mathbf{Y} &\equiv -\mathbf{X} \mathbf{V} \mathbf{W}^{-1} \\ \text{and } \mathbf{Z} &\equiv \mathbf{W}^{-1} + \mathbf{W}^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{V} \mathbf{W}^{-1}, \end{aligned}$$

It is worth noting that the matrix  $\mathbf{Z}$  is not a block-diagonal matrix. However, since the covariance between the fitted values of two different groups is rarely of interest, it suffices

to compute and store the diagonal blocks. Appendix 5.A provides details of the derivations for the streamlined MFVB algorithm for model (5.11).

### 5.3 Numerical evaluation

In this section, we conducted a comprehensive simulation study to evaluate the performance of Algorithm 15 in terms of Bayesian inferential accuracy and computational speed. We simulated 30 independent datasets from model (5.9) and used 25 knots for the overall mean function and 10 knots for the group-specific deviation functions:

$$\begin{aligned}
 y_{ij} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\text{logit}^{-1}\{f(x_{ij}) + g(x_{ij})\}), & (5.14) \\
 \text{where } f(x) &= 10 \phi(x; 1.5, 0.3) + 6 \phi(x; 4, 0.6), \\
 g_i(x) &= \alpha \sin(2\pi x^\beta) \ ; \ x_{ij} \sim \text{Uniform}(0, 1/n) \\
 \text{and } \alpha &\sim N(0, 0.5) \ ; \ \beta \in \{1, 2, 3\}.
 \end{aligned}$$

Each dataset consists of  $m = 50$  groups and the number of observations per group is balanced with  $n_i = n = 50$ . The term  $\phi(x; a, b)$  is an univariate normal density with mean  $a$  and standard deviation  $b$ .

We fitted each replicated dataset using both MFVB approximation via a streamlined version of Algorithm 15 and MCMC sampling. The MFVB algorithm was implemented in the R computing environment (R Development Core Team, 2015) and its stopping criterion is when the relative change in the lower bound  $\underline{p}(\mathbf{y}; q)$  falls below  $10^{-8}$ . The MCMC samples were obtained using the R package `RStan` (Stan Development Team, 2015). In each dataset, MCMC samples of size 10000 were generated, with the first 5000 values of each sample discarding as burn-in and the remaining 5000 values thinning by a factor of 5. Trace plots and the autocorrelation functions for all model parameters were used to assess MCMC convergence.

#### 5.3.1 Assessment of accuracy

For a generic parameter  $\theta$  we assessed the accuracy of the MFVB approximation by comparing the optimal  $q$ -density function  $q^*(\theta)$  with a highly accurate MCMC-based posterior approximation  $p_{\text{MCMC}}(\theta|\mathbf{y})$ . While there are numerous means of measuring accuracy, we recommend working with the a measure that is based on the  $L_1$  distance, also known as the *integrated absolute error* (Faes *et al.*, 2011) of  $q^*(\theta)$ , given by

$$\text{accuracy}\{q^*(\theta)\} \equiv 100 \left( 1 - \frac{1}{2} \int_{-\infty}^{\infty} |q^*(\theta) - p_{\text{MCMC}}(\theta|\mathbf{y})| d\theta \right) \%.$$

This accuracy measure has the attraction of being invariant to monotone transformations on the parameter  $\theta$ . Exact computation of  $p_{\text{MCMC}}(\theta|\mathbf{y})$  is numerically challenging and hence we used binned kernel density estimation with direct plug-in bandwidth selection, as facilitated in the R package `KernSmooth` (Wand and Ripley, 2006).

Figures 5.4 and 5.5 display the approximate posterior density functions and side-by-side boxplots of accuracy scores for the model parameters  $f(Q_k)$ ,  $1 \leq k \leq 3$  and  $g_i(Q_k)$ ,  $k = 2$ , where  $Q_k$  is the  $k$ th sample quintile of the  $x$ s. The boxplots show that the majority of the accuracy scores exceed 80%, with some over 90%. As indicated by Figure 5.4, the MFVB credible intervals have good coverage properties. However, interestingly the MFVB credible intervals tend to be narrower than that of the MCMC samples based on kernel density estimation. The average elapsed time for the MCMC fits is 7.1 hours (standard error 3.0 hours), whilst for the MFVB fits is 1.9 minutes (standard error 4.7 seconds). This corresponds to a speed-up in the order of several hundreds.

### 5.3.2 Assessment of speed

We now turn our attention to quantification of the speed gains afforded by the streamlined MFVB algorithm. The simulation described in Section 5.3 was re-run using MFVB with a combination of  $m$  and  $n$  values and MCMC omitted. All computations were performed on a Mac OS X laptop with a 2.6 GHz Intel Core i7 processor and 4 GBytes of random access memory. Table 5.1 summarises the average (standard error) computing times over 30 runs and highlights the practical and scalability benefits of the streamlined MFVB approach.

Number of groups ( $m$ )	100	500	1000
Within-group size ( $n$ )			
10	65.1 (0.86)	351.6 (5.49)	709.4 (4.42)
50	127.1 (2.57)	826.7 (9.63)	1859.2 (6.15)
100	106.5 (1.51)	929.3 (10.77)	2870.3 (13.62)

Table 5.1: Average (standard error) elapsed of the computing times in seconds for the streamlined MFVB Algorithm15 in the simulation setting described in (5.14).

It is well-established that MCMC can be very slow in situations where complex models are applied to large datasets. For the setting in Table 5.1, we expect the MCMC fitting takes days to weeks to run and therefore a similar timing comparison between MFVB and MCMC is not practical.

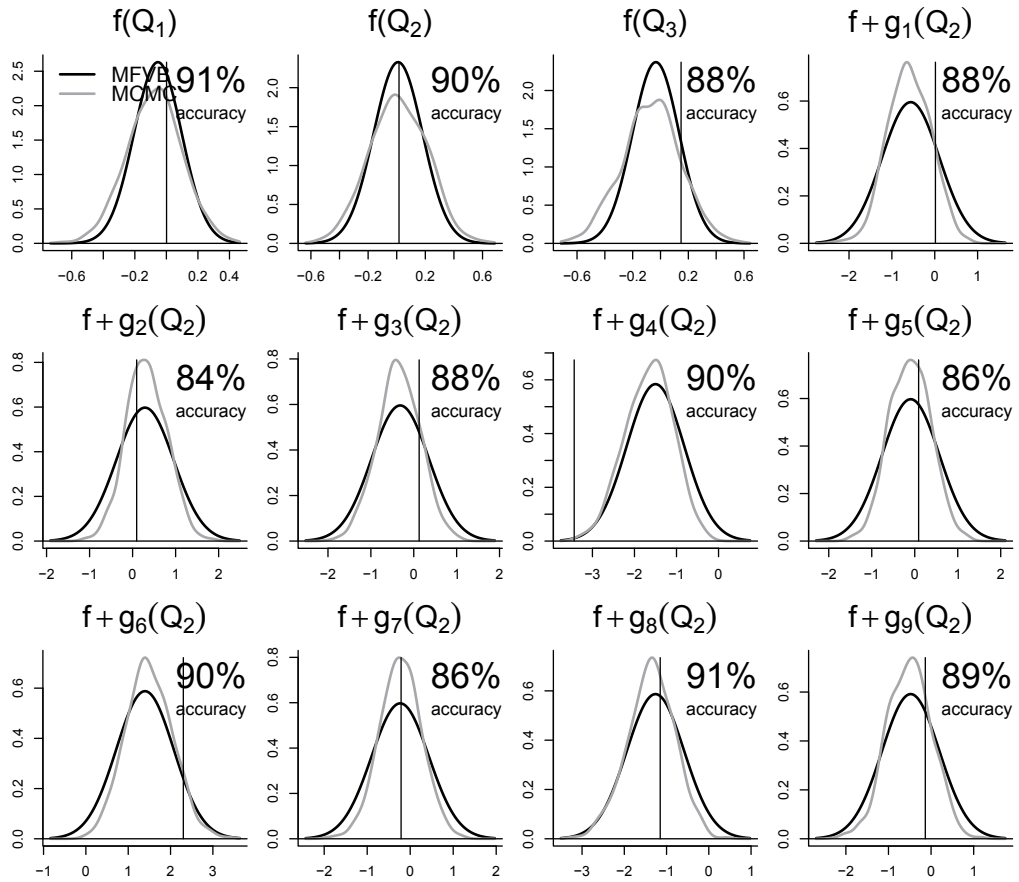


Figure 5.4: The approximate posterior density functions obtained from MFVB and MCMC for a single replication of the simulation study described in the text. Each pair of density function corresponds to a model parameter  $f(Q_k)$ ,  $1 \leq k \leq 3$  and  $g_i(Q_k)$ ,  $k = 2$ , where  $Q_k$  is the  $k$ th sample quintile of the  $x$ s. The vertical lines represent the true parameter values. The accuracy scores on the top right of on each plot show the accuracy of MFVB approximation compared against a MCMC benchmark.

## 5.4 Application to caesarean section data

As described in Section 5.2, we applied our MFVB algorithm to the caesarean section data, with the aim of characterising trends in caesarean section rates for low-risk nulliparous women aged less than 25 years and those aged greater or equal to 25 years in NSW hospitals between 1994 and 2010. A group-specific curve model with a factor-by-curve interaction was fitted to data, allowing for non-linear estimation of time courses as seen in Figure 5.2.

From 1994 to 2010 there were 295 340 low-risk nulliparous women giving birth for the first time in 99 NSW public or private hospitals. Of which 73 795 women (24.5%) aged less than 25 years and 221 545 women (75.0%) aged greater or equal to 25 years. The annual rate of caesarean section among low-risk nulliparous women has increased radically from

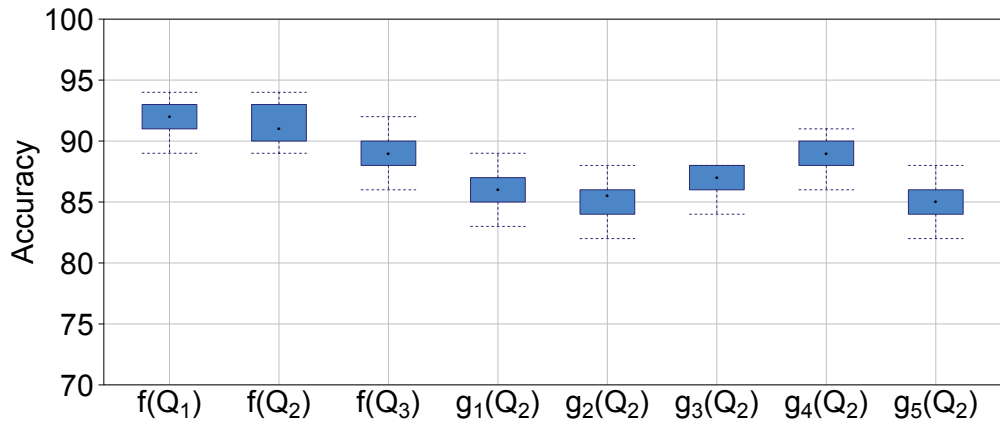


Figure 5.5: Side-by-side boxplots of accuracy scores for the MFVB approximation compared against MCMC over 30 runs. Each boxplot corresponds to a model parameter  $f(Q_k)$ ,  $1 \leq k \leq 3$  and  $g_i(Q_k)$ ,  $k = 2$ , where  $Q_k$  is the  $k$ th sample quintile of the  $x$ s.

12.5% to 24.1% in NSW over this 17-year period; however, rates for women of older age were consistently higher than that of those of younger age (younger group: 10.2% to 16.0%; older group: 13.5% to 26.6%). Figure 5.6 shows the fitted probability function of caesarean section for each hospital between 1994 and 2010. Compared with low-risk women aged less than 25 years, those aged greater or equal to 25 years are more likely to have a caesarean section, rather than a vaginal delivery. The fitted probabilities are seen to vary markedly among hospitals and over time, with some hospitals exhibiting a consistently high probability and others exhibiting an upside down  $U$  shape. It is evident that the hospital-specific probability functions of caesarean section cannot be adequately modelled by linear functions.

Figure 5.7 shows the estimated overall and selected hospital-specific contrast curves defined by (5.12), corresponding to the odds of caesarean section for low-risk nulliparous women aged greater or equal to 25 years, as a function of time, compared with those aged less than 25 years. The pointwise 95% credible sets, shown in the first panel by the grey shaded region, suggest that there are significant differences in trends in caesarean section rates between the two maternal age groups over the year of birth. Low-risk nulliparous women of older age are on average 1.3 times more likely to have a caesarean section compared with those of younger age in the year 1994. This odds ratio increases as the year of birth increases.

To examine hospital differences in odds of caesarean section for low-risk nulliparous women of older age compared with those of younger age, we calculated and plotted the hospital-specific odds ratios for the latest year of birth. In 2010, the statewide odds ratio is 1.74, meaning that, on average, a low-risk nulliparous women of older age giving birth by caesarean section is about 1.74 times more likely than that of a women of younger age.

5.4. APPLICATION TO CAESAREAN SECTION DATA

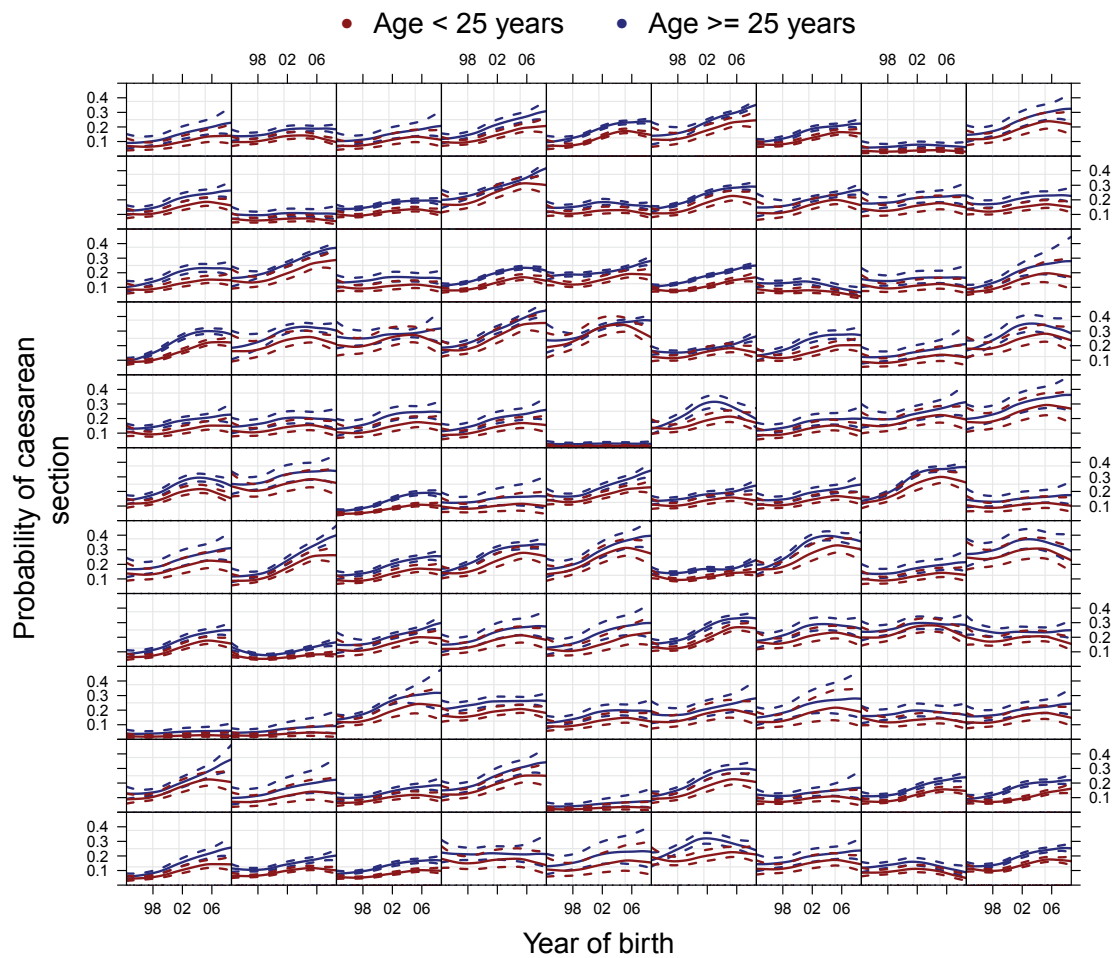


Figure 5.6: The MFVB fitted hospital-specific probability functions of caesarean section for low-risk nulliparous women of aged less than 25 years and those of aged greater or equal to 25 years, as a function of time for each hospital. The dashed curves represent pointwise 95% credible sets. Each panel corresponds to a different hospital.

5.4. APPLICATION TO CAESAREAN SECTION DATA

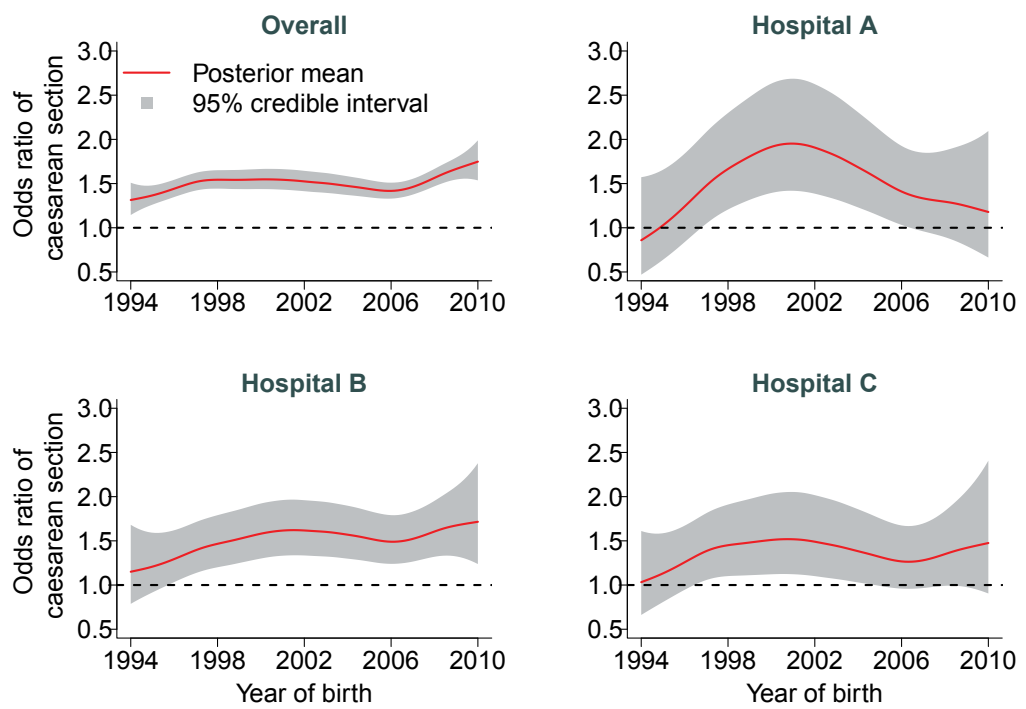


Figure 5.7: The MFVB estimated overall and selected hospital-specific contrast curves defined by (5.12), corresponding to the odds of caesarean section for low-risk nulliparous women aged greater or equal to 25 years, as a function of time, compared with those aged less than 25 years. The shaded regions correspond to pointwise 95% credible sets.

## 5.5. CONCLUDING REMARKS

Year	Odds Ratio	95% Credible Interval
1994	1.31	(1.14, 1.50)
1998	1.42	(1.24, 1.62)
2002	1.52	(1.34, 1.74)
2006	1.63	(1.44, 1.86)
2010	1.74	(1.53, 1.98)

Table 5.2: The MFVB estimated odds ratios of caesarean section (averaged across hospitals) for low-risk nulliparous women aged greater than or equal to 25 years compared with those aged less than 25 years for selected years of birth.

There is a wide variation in odds ratios among hospitals, ranging from 1.21 to 2.24 (Figure 5.8).

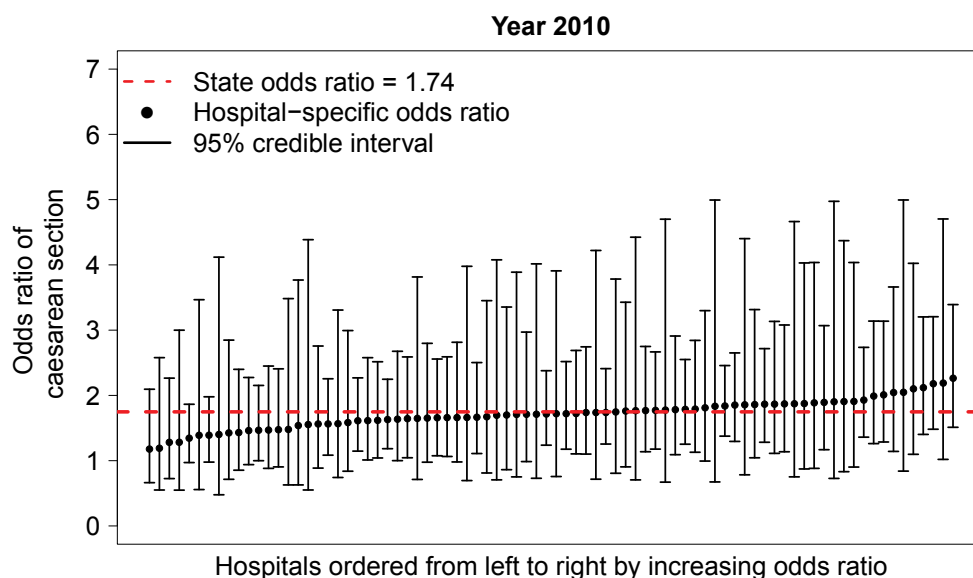


Figure 5.8: The MFVB estimated hospital-specific odds ratios of caesarean section for low-risk nulliparous women of older age compared with those of younger age in the year 2010.

## 5.5 Concluding remarks

We have described MFVB methods for fast approximate inference pertaining to three Bayesian logistic mixed models. Their utility to describe the overall mean and hospital-specific trends in caesarean section rates in NSW between 1994 and 2010 has been demonstrated. Using MCMC methods as a benchmark, we have evaluated the inferential accuracy and computation speed of our streamlined MFVB algorithms.



## 5.5. CONCLUDING REMARKS

---

Compared to MCMC methods which directly generate Markov chain samples from the target posterior distribution, MFVB methods seek an approximating distribution in a factorised form that has minimum Kullback-Leibler distance to the target posterior. Although both methods are an approximation to the true posterior, MCMC methods are theoretically guaranteed to converge asymptotically to the true posterior (but come at great expense), while MFVB methods are known to underestimate the true variance due to the factorisation assumption (Wainwright and Jordan, 2008). Our simulation study shows moderately good to excellent accuracy of the MFVB approximation for all posterior densities of trend parameters. There is some underestimation in the posterior standard deviations, which may also be attributable to the Jaakkola and Jordan’s approach described in (5.7). Despite minor degradation in accuracy, we believe that the MFVB approximations are a strong complement to the standard MCMC methods, as a way to efficiently explore possible models for revealing the complex non-linear temporal evolution of disease rates and identifying regions with unusual patterns at any point in time. Perhaps more practically, the MFVB parameter estimates can be used as starting points for MCMC to speed up convergence. In addition, our simulation results highlight the key advantage of the MFVB methods – computational efficiency. It serves as a great tool for iterative model building/construction. Even with concerns over accuracy in mind, MFVB methods provide an easy way to experiment with a range of models of different specifications, predictors and random effects covariance structures. As illustrated, the computational time for the MFVB algorithms is significantly shorter than that required for MCMC, with answers delivered in minutes rather than hours, although considerations such as computing environment and convergence criterion need to be taken into account to allow a fairer speed comparison. One of the more computationally expensive steps in the MFVB algorithms is the matrix inversion associated with the update of the covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$ , a cost that grows cubically in the number of groups. However, through matrix permutation and block decomposition we derived a streamlined approach that reduces the number of operations to be linear in the number of groups and is memory efficient.

The MFVB methods we consider here are fast and versatile and can be easily extended to more complicated scenarios. For example, the methods allow arbitrary priors for the hyperparameters (Neville *et al.*, 2014) and similar types of models with Gaussian response (Lee and Wand, 2015a; Luts, 2015; Tan and Nott, 2013). Stewart (2014) provides great examples of more elaborate models within the context of social sciences. The ability to handle binary outcomes is particularly useful in health science studies, where binary outcomes are common. For variational inference in logistic regression, we followed on Jaakkola and Jordan (1997) justifying based on constructing a lower bound for the marginal likelihood using convex duality Jaakkola and Jordan (1997) and derived closed-form expressions that approximate the posterior distributions for the parameters,

thus numerical techniques are not required. An extension to our work would be to develop approximations with less stringent conditional independence assumptions to see if the accuracy of MFVB approximations improves. This would be especially useful when we have prior knowledge about the conditional independence structure of the variables. Nevertheless, the  $q$ -density factorisation we investigated in the chapter remains appealing for its simplicity and ease of use.

In this chapter, we presented a flexible estimation of population- and hospital-level trends in caesarean section rates using a penalised spline basis function approach. Apart from being conceptually appealing, the approach allows the resulting models to be couched within a Bayesian mixed model framework for approximate inference. The presented group-specific curve models incorporate both temporal and hospital variabilities into a cohesive modelling framework, enabling one to visually describe the dynamic evolution of outcome rates across hospitals and over time. Our proposed contrast functions are useful in identifying particular regions of the temporal profiles that show significant differences in odds of caesarean section between low-risk nulliparous women of older and younger age. These regions are potential priority targets for public health interventions that aim at reducing temporal and hospital variations in caesarean section rates.

## 5.A Derivation of the mean field variational Bayes algorithm for model (5.11)

The full Bayesian representation for model (5.11) is

$$\begin{aligned}
 \mathbf{y} | \boldsymbol{\beta}, \mathbf{u} &\sim \text{Bernoulli}(\text{logit}^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\beta}^2 \mathbf{I}), \\
 \mathbf{u} | \sigma_{\text{gbl}}^{\text{A}}, \sigma_{\text{gbl}}^{\text{B}}, \boldsymbol{\Sigma}_{\text{R}}, \sigma_{\text{grp}} &\sim N\left(\mathbf{0}, \begin{bmatrix} (\sigma_{\text{gbl}}^{\text{A}})^2 \mathbf{I}_{K_{\text{gbl}}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\sigma_{\text{gbl}}^{\text{B}})^2 \mathbf{I}_{K_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag}_{1 \leq i \leq m} \left( \begin{bmatrix} \boldsymbol{\Sigma}_{\text{R}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sigma_{\text{grp}}^2 \mathbf{I}_{K_{\text{grp}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \sigma_{\text{grp}}^2 \mathbf{I}_{K_{\text{grp}}} \end{bmatrix} \right) \end{bmatrix}\right), \\
 \boldsymbol{\Sigma}_{\text{R}} | a_{\text{R},1}, a_{\text{R},2}, a_{\text{R},3}, a_{\text{R},4} &\sim \text{Inverse-Wishart}(\nu + 3, 2\nu \text{diag}(1/a_{\text{R},1}, \dots, 1/a_{\text{R},4})), \\
 a_{\text{R},1}, a_{\text{R},2}, a_{\text{R},3}, a_{\text{R},4} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, A_{\text{R}}^{-2}\right), \\
 (\sigma_{\text{gbl}}^{\text{A}})^2 | a_{\text{gbl}}^{\text{A}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{gbl}}^{\text{A}}\right), \quad a_{\text{gbl}}^{\text{A}} \sim \text{Inverse-Gamma}\left(\frac{1}{2}, A_{\text{gbl}}^{-2}\right), \\
 (\sigma_{\text{gbl}}^{\text{B}})^2 | a_{\text{gbl}}^{\text{B}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{gbl}}^{\text{B}}\right), \quad a_{\text{gbl}}^{\text{B}} \sim \text{Inverse-Gamma}\left(\frac{1}{2}, A_{\text{gbl}}^{-2}\right), \\
 \sigma_{\text{grp}}^2 | a_{\text{grp}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{grp}}\right), \quad a_{\text{grp}} \sim \text{Inverse-Gamma}\left(\frac{1}{2}, A_{\text{grp}}^{-2}\right),
 \end{aligned}$$

Figure 5.9 shows the directed acyclic graph, as a graph theoretic representation of the Bayesian group-specific curve model conveyed in (5.15). The shaded node corresponds to the observed data vector and the non-shaded nodes correspond to the set of random variables. The advantages to such a graphical representation are that it provides a visualisation of the hierarchical structure of its corresponding Bayesian model, and the graph theoretic results can be used to determine probabilistic relationships between nodes.

Calculations similar to those in Subsection 5.2.1 are used to approximate the full joint posterior density function  $p$  by a product of approximating  $q$ -density functions:

$$\begin{aligned}
 p(\boldsymbol{\beta}, \mathbf{u}, a_{\text{gbl}}^{\text{A}}, a_{\text{gbl}}^{\text{B}}, a_{\text{grp}}^{\text{A}}, a_{\text{grp}}^{\text{B}}, \mathbf{a}_{\text{R}}, \sigma_{\text{gbl}}^{\text{A}}, \sigma_{\text{gbl}}^{\text{B}}, \sigma_{\text{grp}}, \boldsymbol{\Sigma}_{\text{R}}) \\
 \approx q(\boldsymbol{\beta}, \mathbf{u}) q(a_{\text{gbl}}^{\text{A}}) q(a_{\text{gbl}}^{\text{B}}) q(a_{\text{grp}}) q(\mathbf{a}_{\text{R}}) q(\sigma_{\text{gbl}}^{\text{A}}) q(\sigma_{\text{gbl}}^{\text{B}}) q(\sigma_{\text{grp}}) q(\boldsymbol{\Sigma}_{\text{R}}),
 \end{aligned}$$

with the optimal  $q$ -density functions admitting the following forms:

$$\begin{aligned}
 \boldsymbol{\xi} &\leftarrow \sqrt{\text{diagonal}[\mathbf{C}\{\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}^{\text{T}}\} \mathbf{C}^{\text{T}}]}, \\
 q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) &\text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}) \text{ density function,} \\
 q^*(a_{\text{gbl}}^{\text{A}}) &\text{ is the Inverse-Gamma } \left(1, B_{q(a_{\text{gbl}}^{\text{A}})}\right) \text{ density function,} \\
 q^*(a_{\text{gbl}}^{\text{B}}) &\text{ is the Inverse-Gamma } \left(1, B_{q(a_{\text{gbl}}^{\text{B}})}\right) \text{ density function,} \\
 q^*(a_{\text{R},r}) &\text{ is the Inverse-Gamma } \left(1, B_{q(a_{\text{R},r})}\right) \text{ density function, } 1 \leq r \leq 4,
 \end{aligned}$$

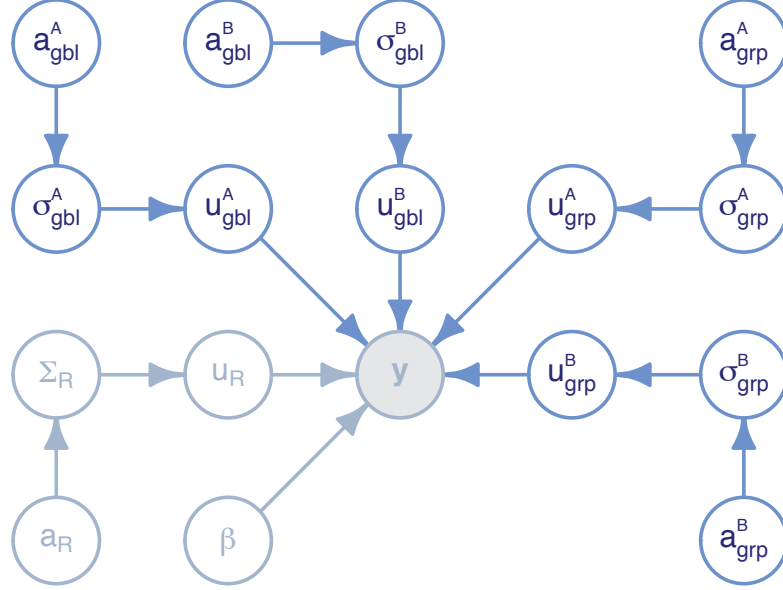


Figure 5.9: Directed Acyclic Graph for the Bernoulli group-specific curve model with contrasting. The shaded node corresponds to the observed data vector and the open nodes correspond to the random or auxiliary variables. The fragments shown in grey are present in the Bernoulli model DAG graph as shown in Figure 2.4.

$q^*((\sigma_{\text{gbl}}^{\text{A}})^2)$  is the Inverse-Gamma  $\left(\frac{1}{2}(K_{\text{gbl}} + 1), B_{q((\sigma_{\text{gbl}}^{\text{A}})^2)}\right)$  density function,  
 $q^*((\sigma_{\text{gbl}}^{\text{B}})^2)$  is the Inverse-Gamma  $\left(\frac{1}{2}(K_{\text{gbl}} + 1), B_{q((\sigma_{\text{gbl}}^{\text{B}})^2)}\right)$  density function,  
 $q^*(\sigma_{\text{grp}}^2)$  is the Inverse-Gamma  $\left(\frac{1}{2}(m K_{\text{grp}} + 1), B_{q(\sigma_{\text{grp}}^2)}\right)$  density function and  
 $q^*(\Sigma_{\text{R}})$  is the Inverse-Wishart  $(\nu + m + 1, \mathbf{B}_{q(\Sigma_{\text{R}})})$  density function.

The parameters  $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}$ ,  $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}$  are the respective mean vector and covariance matrix of  $q^*(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$ . The parameter  $B_{q(\cdot)}$  is the scale parameter of the inverse-Gamma  $q$ -density. Appendix B of Wand and Ormerod (2011) provides a step-by step derivation for the optimal  $q$ -densities of similar type.

### 5.A.1 Mixed model representation

Using the equivalence between penalised splines and mixed models, we now aim to express model (5.11) within the mixed model framework. Define the following linear and non-linear predictor vectors and random effects vectors corresponding to the  $i$ th group to be:

$$\mathbf{X}_i^* \equiv \begin{bmatrix} \mathbf{1} & \mathbf{x}_i \end{bmatrix} ; \quad \mathbf{Z}_{\text{gbl}, i} \equiv \begin{bmatrix} z_{\text{gbl}, 1}(\mathbf{x}_i) & \cdots & z_{\text{gbl}, K_{\text{gbl}}}(\mathbf{x}_i) \end{bmatrix},$$

5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM FOR MODEL (5.11)

---

$$\begin{aligned}
\mathbf{Z}_{\text{grp},i} &\equiv \left[ z_{\text{grp},1}(\mathbf{x}_i) \cdots z_{\text{grp},K_{\text{grp}}}(\mathbf{x}_i) \right], \\
\mathbf{u}_{\text{gbl},i}^{\text{A}} &\equiv \left[ u_{\text{gbl},1}^{\text{A}} \cdots u_{\text{gbl},K_{\text{gbl}}}^{\text{A}} \right]^{\text{T}}; \quad \mathbf{u}_{\text{gbl},i}^{\text{B}} \equiv \left[ u_{\text{gbl},1}^{\text{B}} \cdots u_{\text{gbl},K_{\text{gbl}}}^{\text{B}} \right]^{\text{T}}, \\
\mathbf{u}_{\text{grp},i}^{\text{A}} &\equiv \left[ u_{\text{grp},1}^{\text{A}} \cdots u_{\text{grp},K_{\text{grp}}}^{\text{A}} \right]^{\text{T}}; \quad \mathbf{u}_{\text{grp},i}^{\text{B}} \equiv \left[ u_{\text{grp},1}^{\text{B}} \cdots u_{\text{grp},K_{\text{grp}}}^{\text{B}} \right]^{\text{T}} \\
\text{and} \quad \mathbf{u}_{\text{R},i} &\equiv \left[ U_{0i}^{\text{A}} \quad U_{1i}^{\text{A}} \quad U_{0i}^{\text{B}} \quad U_{1i}^{\text{B}} \right]^{\text{T}},
\end{aligned}$$

where  $\mathbf{x}_i$  is an  $n_i \times 1$  vector containing the  $x_{ij}$ .

In Section 2 of Zhao *et al.* (2006) notation, the compact matrix form of model (5.11) is

$$\boldsymbol{\beta} \equiv \begin{bmatrix} \boldsymbol{\beta}^{\text{R}} \\ \boldsymbol{\beta}^{\text{G}} \end{bmatrix}, \quad \mathbf{X} \equiv [\mathbf{X}^{\text{R}} \quad \mathbf{X}^{\text{G}}], \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}^{\text{R}} \\ \mathbf{u}^{\text{G}} \end{bmatrix} \quad \text{and} \quad \mathbf{Z} \equiv [\mathbf{Z}^{\text{R}} \quad \mathbf{Z}^{\text{G}}],$$

which leads to

$$\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} = \mathbf{X}^{\text{R}}\boldsymbol{\beta}^{\text{R}} + \mathbf{X}^{\text{G}}\boldsymbol{\beta}^{\text{G}} + \mathbf{Z}^{\text{R}}\mathbf{u}^{\text{R}} + \mathbf{Z}^{\text{G}}\mathbf{u}^{\text{G}}.$$

Define  $\mathbf{I}_i^{\text{A}} = \mathbf{1}$  if  $(\mathbf{x}_i, \mathbf{y}_i)$  is of type A and zero otherwise. The matrices  $\mathbf{X}^{\text{R}}$  and  $\mathbf{Z}^{\text{R}}$  are random design matrices corresponding to the fixed effects vector  $\boldsymbol{\beta}^{\text{R}}$  and random group effects vector  $\mathbf{u}^{\text{R}}$  (superscript R) respectively. They are defined as:

$$\begin{aligned}
\mathbf{X}^{\text{R}} &\equiv \emptyset, \quad \mathbf{Z}^{\text{R}} \equiv \text{blockdiag}(\mathbf{X}_i^* \mid \mathbf{I}_i^{\text{A}} \odot \mathbf{Z}_{\text{grp},i} \mid (\mathbf{1} - \mathbf{I}_i^{\text{A}}) \odot \mathbf{Z}_{\text{grp},i}), \\
\boldsymbol{\beta}^{\text{R}} &\equiv \emptyset, \quad \mathbf{u}^{\text{R}} \equiv \left[ (\mathbf{u}_{\text{R},1}^{\text{T}} \quad (\mathbf{u}_{\text{grp},1}^{\text{A}})^{\text{T}} \quad (\mathbf{u}_{\text{grp},1}^{\text{B}})^{\text{T}} \cdots \mathbf{u}_{\text{R},m}^{\text{T}} \quad (\mathbf{u}_{\text{grp},m}^{\text{A}})^{\text{T}} \quad (\mathbf{u}_{\text{grp},m}^{\text{B}})^{\text{T}} \right]^{\text{T}}.
\end{aligned}$$

The matrices  $\mathbf{X}^{\text{G}}$  and  $\mathbf{Z}^{\text{G}}$  are general design matrices corresponding to the fixed effects vector  $\boldsymbol{\beta}^{\text{G}}$  and random spline coefficients vector  $\mathbf{u}^{\text{G}}$  (superscript G) respectively. Typically,  $\mathbf{X}^{\text{G}}$  contains polynomial functions of continuous predictors that are modeled as penalised splines and  $\mathbf{Z}^{\text{G}}$  would then contain random spline basis functions of the same predictors. They are defined as:

$$\begin{aligned}
\boldsymbol{\beta}^{\text{G}} &\equiv \begin{bmatrix} \beta_0^{\text{A}} \\ \beta_1^{\text{A}} \\ \beta_0^{\text{BvsA}} \\ \beta_1^{\text{BvsA}} \end{bmatrix}, \quad \mathbf{X}^{\text{G}} \equiv \begin{bmatrix} \mathbf{X}_i^* & (\mathbf{1} - \mathbf{I}_i^{\text{A}}) \odot \mathbf{X}_i^* \\ \vdots & \vdots \\ \mathbf{X}_m^* & (\mathbf{1} - \mathbf{I}_m^{\text{A}}) \odot \mathbf{X}_m^* \end{bmatrix}, \\
\mathbf{Z}^{\text{G}} &\equiv \begin{bmatrix} \mathbf{I}_i^{\text{A}} \odot \mathbf{Z}_{\text{gbl},i} & (\mathbf{1} - \mathbf{I}_i^{\text{A}}) \odot \mathbf{Z}_{\text{gbl},i} \\ \vdots & \vdots \\ \mathbf{I}_m^{\text{A}} \odot \mathbf{Z}_{\text{gbl},m} & (\mathbf{1} - \mathbf{I}_m^{\text{A}}) \odot \mathbf{Z}_{\text{gbl},m} \end{bmatrix}, \\
\text{and} \quad \mathbf{u}^{\text{G}} &\equiv \left[ (\mathbf{u}_{\text{gbl},1}^{\text{A}})^{\text{T}} \cdots (\mathbf{u}_{\text{gbl},m}^{\text{A}})^{\text{T}} \quad (\mathbf{u}_{\text{gbl},1}^{\text{B}})^{\text{T}} \cdots (\mathbf{u}_{\text{gbl},m}^{\text{B}})^{\text{T}} \right]^{\text{T}},
\end{aligned}$$

### 5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM FOR MODEL (5.11)

---

where  $\odot$  denotes the element-wise product of matrices.

Building on Lee and Wand (2015a), we derive a streamlined iterative scheme for obtaining optimal  $q$ -density moments for all model parameters of the group-specific curve model with a factor-by-curve interaction. Our presentation of the streamlined MFVB algorithm benefits from the following notation, where  $\mathbf{y}$ ,  $\mathbf{X}^G$ ,  $\mathbf{Z}^R$  and  $\mathbf{Z}^G$  are partitioned row-wise corresponding to the  $i$ th group:

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X}^G = \begin{bmatrix} \mathbf{X}_1^G \\ \vdots \\ \mathbf{X}_m^G \end{bmatrix}, \quad \mathbf{Z}^R = \begin{bmatrix} \mathbf{Z}_1^R \\ \vdots \\ \mathbf{Z}_m^R \end{bmatrix} \quad \text{and} \quad \mathbf{Z}^G = \begin{bmatrix} \mathbf{Z}_1^G \\ \vdots \\ \mathbf{Z}_m^G \end{bmatrix}.$$

Here  $\mathbf{y}_i \equiv [y_{i1} \cdots y_{in_i}]^T$  denotes the  $n_i \times 1$  vector of responses for the  $i$ th group. The matrices  $\mathbf{X}_i^G$ ,  $\mathbf{Z}_i^R$  and  $\mathbf{Z}_i^G$  are defined in the same fashion. In addition, it is useful to define

$$\mathbf{C}^G \equiv [\mathbf{X} \ \mathbf{Z}^G] \quad \text{and} \quad \mathbf{C}^R \equiv [\mathbf{X}^R \ \mathbf{Z}^G].$$

#### 5.A.2 Streamlined mean field variational Bayes algorithm

We are now ready to present the streamlined algorithm for fast MFVB approximate fitting and inference for model (5.11).

#### 5.A.3 Fitting the group-specific curve model (5.11) in Rstan

Stan is a probabilistic programming language, written in C++, for implementing full Bayesian statistical inference through Hamiltonian Monte Carlo No U-turn sampling, a form of MCMC sampling. We now provide the `Rstan` code for fitting of model (5.11).

Specify the data: the number of observation,  $\sum_{i=1}^m n_i$ ; the number of hospitals,  $m$ ; the hospital identification number, `idnum`; the response vector,  $\mathbf{y}$ ; the respective fixed and random effects design matrices,  $\mathbf{X}$  and  $\mathbf{Z}$ ; and the hyperparameters,  $\sigma_\beta^2$ ,  $A_R$ ,  $A_{\text{gbl}}$  and  $A_{\text{grp}}$ . Data are labelled as integer or real and can be vectors if dimensions are specified. Data can also be constrained; for example all hyperparameters must be positive.

```
grpSpecCurveModel <-
'data
{
  int<lower=1> numObs;           int<lower=1> m;
  int<lower=1> idnum[numObs];   int<lower=1> ncXR;
  int<lower=1> ncZ;             int<lower=1> ncZgrp;
  real<lower=0> sigmaBeta;      real<lower=0> AR;
  real<lower=0> AuGbl;          real<lower=0> AuGrp;
  matrix[numObs,ncXR] Xbase;   matrix[numObs,ncXR] XTypA;
```

5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM FOR MODEL (5.11)

---

**Initialise:**  $\nu = 4$ ,  $\mu_{q(1/(\sigma_{\text{gbl}}^{\text{A}})^2)} > 0$ ,  $\mu_{q(1/(\sigma_{\text{gbl}}^{\text{B}})^2)} > 0$ ,  $\mu_{q(1/a_{\text{gbl}}^{\text{A}})} > 0$ ,  $\mu_{q(1/a_{\text{gbl}}^{\text{B}})} > 0$ ,  $\mu_{q(1/\sigma_{\text{grp}}^2)} > 0$ ,  $\mu_{q(1/a_{\text{grp}})} > 0$ ,  $\mu_{q(1/a_{\text{R},r})} > 0$ ,  $\mathbf{M}_{q(\Sigma_{\text{R}}^{-1})}$ , a  $q^{\text{R}} \times q^{\text{R}}$  positive definite matrix and  $\boldsymbol{\xi}$ , a  $(\sum_{i=1}^m n_i) \times 1$  vector of positive entries.

**Cycle through updates:**

$$\text{Define: } \mathbf{M}_{q(\Omega)} \equiv \begin{bmatrix} \sigma_{\beta}^{-2} \mathbf{I}_p & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/(\sigma_{\text{gbl}}^{\text{A}})^2)} \mathbf{I}_{K_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mu_{q(1/(\sigma_{\text{gbl}}^{\text{B}})^2)} \mathbf{I}_{K_{\text{gbl}}} \end{bmatrix}$$

$$\mathbf{S} \leftarrow \mathbf{0} \ ; \ \mathbf{s} \leftarrow \mathbf{0}$$

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

$$\begin{aligned} \mathbf{G}_i &\leftarrow 2(\mathbf{C}_i^{\text{G}})^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_i)\} \mathbf{C}_i^{\text{R}} \\ \mathbf{H}_i &\leftarrow \left\{ 2(\mathbf{C}_i^{\text{R}})^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_i)\} \mathbf{C}_i^{\text{R}} + \begin{bmatrix} \mathbf{M}_{q(\Sigma_{\text{R}}^{-1})} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{2K_{\text{grp}}} \end{bmatrix} \right\}^{-1} \\ \mathbf{S} &\leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^{\text{T}} \ ; \ \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{C}_i^{\text{R}})^{\text{T}} (\mathbf{y}_i - \frac{1}{2} \mathbf{1}) \end{aligned}$$

**Update multivariate Normal  $q^*(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})$  parameters:**

$$\Sigma_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \leftarrow \{2(\mathbf{C}^{\text{G}})^{\text{T}} \text{diag}\{\lambda(\boldsymbol{\xi}_i)\} \mathbf{C}^{\text{G}} + \mathbf{M}_{q(\Omega)} - \mathbf{S}\}^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \leftarrow \Sigma_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \{(\mathbf{C}^{\text{G}})^{\text{T}} (\mathbf{y} - \frac{1}{2} \mathbf{1}) - \mathbf{s}\}$$

**Update multivariate Normal  $q^*(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})$  parameters:**

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

$$\begin{aligned} \Sigma_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})} &\leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^{\text{T}} \Sigma_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \mathbf{G}_i \mathbf{H}_i \\ \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})} &\leftarrow \mathbf{H}_i \left\{ (\mathbf{X}_i^{\text{R}})^{\text{T}} \mathbf{y}_i - \mathbf{G}_i^{\text{T}} \boldsymbol{\mu}_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \right\} \end{aligned}$$

**Update  $\boldsymbol{\xi}$  parameters:**

$$\boldsymbol{\xi}^2 \leftarrow \text{diagonal}\{ \mathbf{C}^{\text{G}} (\Sigma_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} + \boldsymbol{\mu}_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \boldsymbol{\mu}_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})}^{\text{T}}) (\mathbf{C}^{\text{G}})^{\text{T}} \}$$

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

$$\begin{aligned} \boldsymbol{\xi}_i^2 &\leftarrow \boldsymbol{\xi}_i^2 + 2 \text{diagonal}\{ \mathbf{C}^{\text{G}} (-\Sigma_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \mathbf{G}_i \mathbf{H}_i \\ &\quad + \boldsymbol{\mu}_{q(\beta, \mathbf{u}_{\text{gbl}}^{\text{A}}, \mathbf{u}_{\text{gbl}}^{\text{B}})} \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})}^{\text{T}}) (\mathbf{C}_i^{\text{R}})^{\text{T}} \} \\ \boldsymbol{\xi}_i^2 &\leftarrow \boldsymbol{\xi}_i^2 + \text{diagonal}\{ (\mathbf{C}_i^{\text{R}} (\Sigma_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})} \\ &\quad + \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})} \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i}^{\text{A}}, \mathbf{u}_{\text{grp},i}^{\text{B}})}^{\text{T}}) (\mathbf{C}_i^{\text{R}})^{\text{T}} \} \end{aligned}$$

5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM FOR MODEL (5.11)

---

**Update inverse-Gamma  $q^*(\sigma_{\text{gbl}}^A)$  and inverse-Gamma  $q^*(a_{\text{gbl}}^A)$  parameters:**

$$B_{q((\sigma_{\text{gbl}}^A)^2)} \leftarrow \mu_{q(1/a_{\text{gbl}}^A)} + \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}_{\text{gbl}}^A)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_{\text{gbl}}^A)}) \right\}$$

$$\mu_{q(1/(\sigma_{\text{gbl}}^A)^2)} \leftarrow \frac{1}{2}(K_{\text{gbl}} + 1)/B_{q((\sigma_{\text{gbl}}^A)^2)} \quad ; \quad \mu_{q(1/a_{\text{gbl}}^A)} \leftarrow 1/(\mu_{q(1/(\sigma_{\text{gbl}}^A)^2)} + A_{\text{gbl}}^{-2})$$

**Update inverse-Gamma  $q^*(\sigma_{\text{gbl}}^B)$  and inverse-Gamma  $q^*(a_{\text{gbl}}^B)$  parameters:**

$$B_{q((\sigma_{\text{gbl}}^B)^2)} \leftarrow \mu_{q(1/a_{\text{gbl}}^B)} + \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}_{\text{gbl}}^B)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_{\text{gbl}}^B)}) \right\}$$

$$\mu_{q(1/(\sigma_{\text{gbl}}^B)^2)} \leftarrow \frac{1}{2}(K_{\text{gbl}} + 1)/B_{q((\sigma_{\text{gbl}}^B)^2)} \quad ; \quad 1/(\mu_{q(1/(\sigma_{\text{gbl}}^B)^2)} + A_{\text{gbl}}^{-2})$$

**Update inverse-Gamma  $q^*(\sigma_{\text{grp}})$  and inverse-Gamma  $q^*(a_{\text{grp}})$  parameters:**

$$B_{q(\sigma_{\text{grp}}^2)} \leftarrow \mu_{q(1/a_{\text{grp}})} + \frac{1}{2} \sum_{i=1}^m \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}_{\text{grp},i}^A, \mathbf{u}_{\text{grp},i}^B)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_{\text{grp},i}^A, \mathbf{u}_{\text{grp},i}^B)}) \right\}$$

$$\mu_{q(1/\sigma_{\text{grp}}^2)} \leftarrow \frac{1}{2}(m K_{\text{grp}} + 1)/B_{q(\sigma_{\text{grp}}^2)}$$

$$B_{q(a_{\text{grp}})} \leftarrow \mu_{q(1/\sigma_{\text{grp}}^2)} + A_{\text{grp}}^{-2} \quad ; \quad \mu_{q(1/a_{\text{grp}})} \leftarrow 1/B_{q(a_{\text{grp}})}$$

**Update inverse-Gamma  $q^*(a_{\text{R},r})$  and inverse-Gamma  $q^*(\boldsymbol{\Sigma}_{\text{R}})$  parameters:**

For  $r = 1, \dots, 4$ :

$$B_{q(a_{\text{R},r})} \leftarrow \nu(M_{q(\boldsymbol{\Sigma}_{\text{R}}^{-1})})_{rr} + A_{\text{R}r}^{-2} \quad ; \quad \mu_{q(1/a_{\text{R},r})} \leftarrow \frac{1}{2}(\nu + 3)/B_{q(a_{\text{R},r})}$$

$$M_{q(\boldsymbol{\Sigma}_{\text{R}}^{-1})} \leftarrow (\nu + m + 1) \mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{R}})}^{-1}$$

$$\mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{R}})} \leftarrow \sum_{i=1}^m (\boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i})} \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i})}^{\text{T}} + \boldsymbol{\Sigma}_{q(\mathbf{u}_{\text{R},i})}) + 2\nu \text{diag}(\mu_{q(1/a_{\text{R},1})}, \dots, \mu_{q(1/a_{\text{R},4})})$$

until the increase in  $p(\mathbf{y}; q)$  is negligible.

**Update the  $q$ -density covariance matrix for the  $i$ th group:**

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

$$\begin{aligned} \boldsymbol{\Lambda}_{q(\boldsymbol{\beta}, \mathbf{u}_{\text{gbl}}^A, \mathbf{u}_{\text{gbl}}^B, \mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i})} &\equiv E_q \left[ \left\{ \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}_{\text{gbl}}^A \\ \mathbf{u}_{\text{gbl}}^B \end{bmatrix} - \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}_{\text{gbl}}^A, \mathbf{u}_{\text{gbl}}^B)} \right\} \left\{ \begin{bmatrix} \mathbf{u}_{\text{R},i} \\ \mathbf{u}_{\text{grp},i} \end{bmatrix} - \boldsymbol{\mu}_{q(\mathbf{u}_{\text{R},i}, \mathbf{u}_{\text{grp},i})} \right\}^{\text{T}} \right] \\ &\leftarrow -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}_{\text{gbl}}^A, \mathbf{u}_{\text{gbl}}^B)} \mathbf{G}_i \mathbf{H}_i \end{aligned}$$

---

Algorithm 16: Mean field variational Bayes algorithm for the Bernoulli group-specific curve model with contrasting.



## 5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM FOR MODEL (5.11)

---

```

matrix[numObs,ncXR] XTypB;      int<lower=0,upper=1> y[numObs];
real x[numObs];                matrix[numObs,ncZ] ZTypA;
matrix[numObs,ncZ] ZTypB;      real ZgrpTypA[numObs,ncZgrp];
real ZgrpTypB[numObs,ncZgrp];  vector[2*ncXR] zeroVec;
}

```

Specify the model parameters: the unknown to be estimated in the model fit. These are the fixed effects vector,  $\beta$ ; the random effects vectors,  $\mathbf{u}_{\text{gbl}}$ ,  $\mathbf{u}_{\text{grp}}$  and  $\mathbf{u}_{\text{R}}$ ; the covariance vectors and matrix,  $\sigma_{\text{gbl}}^2$ ,  $\sigma_{\text{grp}}^2$  and  $\Sigma_{\text{R}}$ . In addition, we parametrize  $\mu \equiv \mathbf{X}\beta + \mathbf{Z}\mathbf{u}$  to be a transformation parameter in order to ensure the sampler runs more efficiently.

```

parameters
{
  vector[ncXR] beta;           vector[ncXR] betaTypB;
  vector[ncZ] uGblTypA;       vector[ncZ] uGblTypB;
  vector[2*ncXR] uR[m];       real uGrpTypA[m,ncZgrp];
  real uGrpTypB[m,ncZgrp];    cov_matrix[2*ncXR] SigmaR;
  vector[2*ncXR] aR;          real<lower=0> sigmauGblTypA;
  real<lower=0> sigmauGblTypB; real<lower=0> sigmauGrp;
}
transformed parameters
{
  vector[numObs] fmean;        vector[numObs] fullMean;
  fmean <- (Xbase*beta + XTypB*betaTypB
            + ZTypA*uGblTypA + ZTypB*uGblTypB);
  for (iA11 in 1:numObs)
    fullMean[iA11] <- (fmean[iA11]
                      + uR[idnum[iA11],3]*XTypA[iA11,1]
                      + uR[idnum[iA11],4]*XTypA[iA11,2]
                      + uR[idnum[iA11],1]*XTypB[iA11,1]
                      + uR[idnum[iA11],2]*XTypB[iA11,2]
                      + dot_product(uGrpTypA[idnum[iA11]],ZgrpTypA[iA11])
                      + dot_product(uGrpTypB[idnum[iA11]],ZgrpTypB[iA11]));
}

```

Specify the model statement: The `Bernoulli` specifies that the response vector  $\mathbf{y}$  has a bernoulli distribution with mean  $\text{logit}^{-1}(\mu)$ , where the mean is specified to be the sum of variables for the overall mean, maternal age deviation from that overall mean and hospital deviations.

```

model
{
  matrix[2*ncXR,2*ncXR] rateForWish;

  y ~ bernoulli_logit(fullMean);

  for (i in 1:m)
    uR[i] ~ multi_normal(zeroVec,SigmaR);

  uGblTypA ~ normal(0,sigmauGblTypA);
  uGblTypB ~ normal(0,sigmauGblTypB);
}

```

5.A. DERIVATION OF THE MEAN FIELD VARIATIONAL BAYES ALGORITHM  
FOR MODEL (5.11)

---

```
for (i in 1:m)
{
  for (k in 1:ncZgrp)
  {
    uGrpTypA[i,k] ~ normal(0,sigmauGrp);
    uGrpTypB[i,k] ~ normal(0,sigmauGrp);
  }
}

rateForWish <- rep_matrix(0,4,4);
for (r in 1:(2*ncXR))
{
  aR[r] ~ inv_gamma(0.5,pow(AR,-2));
  rateForWish[r,r] <- 4/aR[r];
}
SigmaR ~ inv_wishart(5,rateForWish);

beta ~ normal(0,sigmaBeta);
betaTypB ~ normal(0,sigmaBeta);
sigmauGblTypA ~ cauchy(0,AuGbl);
sigmauGblTypB ~ cauchy(0,AuGbl);
sigmauGrp ~ cauchy(0,AuGrp);
},'
```

## Chapter 6

# Alternative Approach Based on Variational Message Passing

*Far better an approximate answer to the right question, which is often vague, than the exact answer to the wrong question, which can always be made precise.*

John Tukey

### 6.1 Introduction

In the previous chapters, we presented the full description and implementation of the state-of-the-art MFVB algorithms in longitudinal and multilevel data settings and showed how MFVB approximations can be flexibly scaled to other classes of statistical models including, for example, higher-level random effects, missing data and measurement error models, and group-specific curve models with contrasting. In this chapter, we introduce an alternative approach, which produces the same posterior approximation as MFVB but with a different algebraic system, known as the *variational message passing* (VMP) (e.g. Winn and Bishop, 2005).

We revisit MFVB for Bayesian semiparametric mixed regression but instead work within the VMP framework. Although the MFVB and VMP approaches each lead to different iterative variational algorithms, they converge to the identical approximating posterior density functions in a wide range of models under the same mean-field restrictions. In essence, the approach is guided by the notion of *message passing*, a general principle that is well known in software engineering for efficient computing within distributed systems, but not in mainstream statistics. As explained in Minka (2005), the

parameter update expressions in the MFVB algorithm can be alternatively expressed as *messages* passed between nodes on a suitable factor graph. Factor graphs, earlier introduced in Frey (1998), is a graphical concept tailored to the VMP approach. Often, the VMP messages are model parameters from the exponential family density functions and are generally summarised in natural canonical forms. As pointed out by Wand (2015), the natural parametrisation of VMP messages leads to greater efficiency gains from using VMP over MFVB in terms of modularisation, since the extension to arbitrarily large models is much more straightforward via the introduction of factor graph fragments. The current version of Infer.NET supports VMP fitting of Bayesian hierarchical models, but is only applicable to certain classes of standard models. For the semiparametric mixed regression scenarios, self-implementation of VMP algorithms is our only option.

Contemporary literature on VMP are mainly contributed by Winn and Bishop (2005), Minka (2005) and Minka *et al.* (2009). However, each article offers a slightly different notation system and description of the VMP scheme. We use the same notation convention as described in Minka (2005). The thrust of this chapter is to provide a unified and notationally friendly framework for the general VMP approach regarding the Gaussian response examples considered in Sections 2.2.3 and 3.7.

Sections 6.2 and 6.3 provide relevant background materials relevant to VMP. We provide a general notationally friendly recipe of the VMP approach in Section 6.5. In Section 6.6 we use a simple Bayesian mixed model to convey the main ideas of VMP and then describe the ease of extension to larger statistical models in Section 6.7. Concluding remarks are given in Section 6.8.

## 6.2 Natural canonical forms for exponential family densities

The natural canonical form of any univariate exponential family distribution in a scalar random variable  $x$  can be written as

$$p(x) \equiv \exp\{\mathbf{T}(x)^\top \boldsymbol{\eta} - A(\boldsymbol{\eta})\} h(x), \quad (6.1)$$

where  $\mathbf{T}(x)$  is the *natural statistic*,  $\boldsymbol{\eta}$  is the *natural parameter*,  $A(\boldsymbol{\eta})$  is the *log-partition function* and  $h(x)$  is the *base measure*. In general, the natural statistic is not unique and is commonly expressed  $\mathbf{T}(x)$  as the simplest possible algebraic form given  $p(x)$ . For example, the Inverse-Gamma( $A, B$ ) density function with shape parameter  $A > 0$  and scale parameter  $B > 0$  has the general form

$$p(x; A, B) = B^A \Gamma(A)^{-1} x^{-A-1} \exp(-B/x), \quad x > 0, \quad (6.2)$$

$$= \exp \left\{ \begin{bmatrix} \log(x) \\ 1/x \end{bmatrix}^\top \begin{bmatrix} -A - 1 \\ -B \end{bmatrix} + A \log(B) - \log \Gamma(A) \right\}.$$

Straightforward algebra shows that (6.2) can be expressed in the form of (6.1) with

$$\mathbf{T}(x) = \begin{bmatrix} \log(x) \\ 1/x \end{bmatrix}, \quad \boldsymbol{\eta} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} -A - 1 \\ -B \end{bmatrix} \quad \text{and} \quad h(x) = 1.$$

The log-partition function is

$$A(\boldsymbol{\eta}) = \log \Gamma(-\eta_1 - 1) - (-\eta_1 - 1) \log(-\eta_2).$$

Table 6.1 lists the natural statistics, natural parameter vectors and natural canonical forms of common exponential family densities.

Probability distribution	Natural statistic	Natural parameter vector	Natural canonical form
Bernoulli( $p$ )	$x$	$\log\{p(1-p)\}$	Bernoulli <sub>nat</sub> ( $\eta$ )
$N(\mu, \sigma^2)$	$\begin{bmatrix} x \\ x^2 \end{bmatrix}$	$\begin{bmatrix} \mu/\sigma^2 \\ -1/(2\sigma^2) \end{bmatrix}$	$N_{\text{nat}}(\eta_1, \eta_2)$
Inverse-Gamma( $A, B$ )	$\begin{bmatrix} \log(x) \\ 1/x \end{bmatrix}$	$\begin{bmatrix} -A - 1 \\ B \end{bmatrix}$	Inverse-Gamma <sub>nat</sub> ( $\eta_1, \eta_2$ )
$N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	$\begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}$	$\begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}$	$N_{\text{nat}}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2,\text{vec}})$
Inverse-Wishart( $A, \mathbf{B}$ )	$\begin{bmatrix} \log(\mathbf{X}) \\ \text{vec}(\mathbf{X}^{-1}) \end{bmatrix}$	$\begin{bmatrix} \frac{1}{2}(A + d + 1) \\ -\frac{1}{2}\text{vec}(\mathbf{B}) \end{bmatrix}$	Inverse-Wishart <sub>nat</sub> ( $\eta_1, \boldsymbol{\eta}_2$ )

Table 6.1: Expressions for natural statistics, natural parameter vectors and natural canonical forms of common exponential family densities.

### 6.3 Primitive integrals and function definitions

The following primitives and function definitions follow immediately from the results involving moments of random variables which take on distributions in Table 6.1. These moments aid in the derivation of the VMP algorithms described in the later sections.

We define

$p_{N_{\text{nat,vec}}}\left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}\right)$  to be the  $N_{\text{nat,vec}}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2,\text{vec}})$  density function in  $\mathbf{x}$ ,

$p_{\text{IG}_{\text{nat}}}\left(x, \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}\right)$  to be the Inverse-Gamma<sub>nat</sub>( $\eta_1, \eta_2$ ) density function in  $x$ ,

$p_{\text{IW}_{\text{nat}}}\left(\mathbf{X}, \begin{bmatrix} \eta_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}\right)$  to be the Inverse-Wishart<sub>nat</sub>( $\eta_1, \boldsymbol{\eta}_2$ ) density function in  $\mathbf{X}$ .

We now present the corresponding algorithmic primitives and function definitions used throughout this chapter:

**Primitive 6.3.1.** If  $\mathbf{x}$  is a  $d \times 1$  random vector,

$$\begin{aligned} \mathcal{F}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}\right) &\equiv \int_{\mathbb{R}^d} \mathbf{x}^\top \mathbf{x} p_{N_{\text{nat,vec}}}\left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}\right) d\mathbf{x} \\ &= \frac{1}{4} \text{tr}\left(\{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1} \left[\boldsymbol{\eta}_1 \boldsymbol{\eta}_1^\top \{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1} - \mathbf{I}\right]\right). \end{aligned}$$

**Primitive 6.3.2.** If  $\mathbf{x}$  is a  $d \times 1$  random vector,  $\mathbf{b}$  a  $p \times 1$  vector and  $\mathbf{C}$  a  $p \times q$  matrix,

$$\begin{aligned} \mathcal{H}\left(\begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}; \mathbf{b}, \mathbf{C}\right) &\equiv \int_{\mathbb{R}^d} (\mathbf{b} - \mathbf{C}\mathbf{x})(\mathbf{b} - \mathbf{C}\mathbf{x})^\top p_{N_{\text{nat,vec}}}\left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}\right) d\mathbf{x} \\ &= \mathbf{b}\mathbf{b}^\top + \{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1} \left[\mathbf{b}\boldsymbol{\eta}_1^\top \mathbf{C}^\top + \frac{1}{4} \mathbf{C}\boldsymbol{\eta}_1 \boldsymbol{\eta}_1^\top \mathbf{C}^\top \{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1} - \frac{1}{2} \mathbf{C}\mathbf{C}^\top\right]. \end{aligned}$$

**Primitive 6.3.3.** If  $\mathbf{X}$  is a  $d \times d$  random vector,

$$\begin{aligned} \mathcal{G}\left(\begin{bmatrix} \eta_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}\right) &\equiv \int_S \mathbf{X}^{-1} p_{\text{IW}_{\text{nat}}}\left(\mathbf{X}; \begin{bmatrix} \eta_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}\right) d\mathbf{x} \\ &= \{\eta_1 + \frac{1}{2}(d+1)\} \{\text{vec}^{-1}(\boldsymbol{\eta}_2)\}^{-1}, \end{aligned}$$

where  $S$  is the set of all symmetric, positive definite  $d \times d$  matrices.

**Primitive 6.3.4.** If  $x$  is a scalar random variable,

$$\mathcal{I} \left( \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \right) \equiv \int_0^\infty \frac{1}{x} p_{\text{IG}_{\text{nat}}} \left( x, \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \right) dx = \frac{\eta_1 + 1}{\eta_2}.$$

## 6.4 Factor graphs

A *factor graph* is a graphical representation of the mathematical relation between arguments and local factors of a real-valued function. Consider, for example, the function  $g(x_1, x_2, x_3, x_4, x_5)$  of five variables defined on  $\mathbb{R}^5$  as follows:

$$g(x_1, x_2, x_3, x_4, x_5) \equiv f_A(x_1) f_B(x_2, x_3) f_C(x_3, x_4, x_5) f_D(x_3, x_4) f_E(x_5). \quad (6.3)$$

Figure 6.1 is a factor graph corresponding to  $g$ . The circular nodes correspond to the arguments of  $g$  and the square nodes correspond to the factors in (6.3). Edges are drawn between each factor node and arguments of that factor. In general, factor graphs of functions are non-unique since, for example,  $f_A$  and  $f_B$  could be aggregated into a single factor resulting in a different factor graph. For the remainder of this chapter, we refer a circular node as a *stochastic node* (i.e. random variables, random vectors or random matrices) and a square node as a *factor*. We use the word *node* to describe either a stochastic node or a factor.

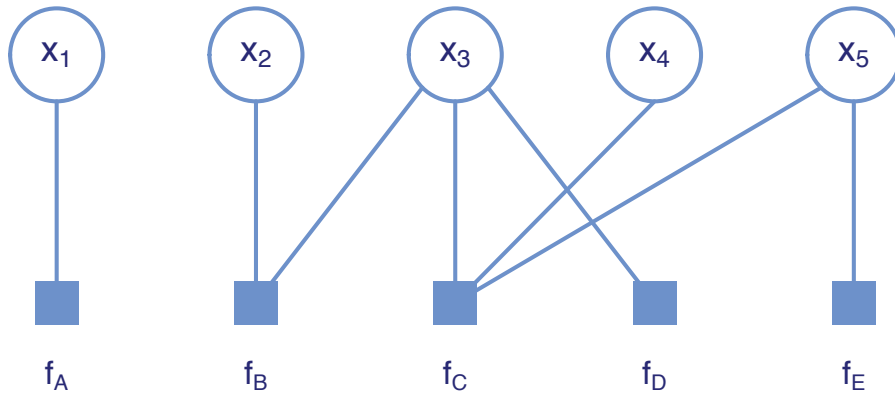


Figure 6.1: A factor graph corresponding to the function  $g(x_1, x_2, x_3, x_4, x_5)$  defined by (6.3).

## 6.5 Variational message passing

Consider a generic Bayesian model with observed data vector  $\mathbf{y}$  and parameter vector  $\boldsymbol{\theta}$  that is continuous over the parameter space  $\Theta$ . A MFVB approximation to the posterior density function  $p(\boldsymbol{\theta}|\mathbf{y})$  is

$$p(\boldsymbol{\theta}|\mathbf{y}) \approx q^*(\boldsymbol{\theta}),$$

where  $q^*(\boldsymbol{\theta})$  is the minimiser of the Kullback-Leibler divergence of  $p(\boldsymbol{\theta}|\mathbf{y})$  from  $q(\boldsymbol{\theta})$  subject to the product density restriction for a subset  $S$  of  $\{1, \dots, M\}$

$$q(\boldsymbol{\theta}) = \prod_{i=1}^M q_i(\boldsymbol{\theta}_S), \quad \boldsymbol{\theta}_S \equiv \{\boldsymbol{\theta}_j : j \in S\}. \quad (6.4)$$

The restriction (6.4) means that the joint density function of  $\boldsymbol{\theta}$  and  $\mathbf{y}$  can be expressed as

$$p(\boldsymbol{\theta}, \mathbf{y}) = \prod_{j=1}^N f_j(\boldsymbol{\theta}_{S_j}) \text{ for subsets } S_j \text{ of } \{1, \dots, M\} \text{ and factors } f_j, 1 \leq j \leq N. \quad (6.5)$$

Variational message passing is a message-passing algorithmic implementation of the mean-field approximation. The VMP iterative updates are expressed in terms of updating messages passed between nodes on the factor graph. For example, if  $p(\boldsymbol{\theta}, \mathbf{y})$  is a directed acyclic graphical model with nodes  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M$  and  $\mathbf{y}$  then

$$p(\boldsymbol{\theta}, \mathbf{y}) = \left\{ \prod_{j=1}^M p(\boldsymbol{\theta}_j | \text{parents of } \boldsymbol{\theta}_j) \right\} p(\mathbf{y} | \text{parents of } \mathbf{y}) \quad (6.6)$$

is an example of (6.5) with the factors  $f_j$  corresponding to the density function of  $\boldsymbol{\theta}_j$  conditional on its parent nodes and the factor  $f_{M+1}$  corresponding to the likelihood. Each factor is a function of the subset of (6.4) representing the parent-child relationships in the directed acyclic graph infrastructure. Figure 6.2 is an example of the factor graph that represents the full factorisation of  $p(\boldsymbol{\theta}|\mathbf{y})$  with respect to (6.6) with  $M = 9$  and  $N = 10$ . It becomes clear that in factor graphs the conditional independence properties of statistical models are represented by the squared factor nodes in replacement of the directed arrows in directed acyclic graphs, henceforth we introduce a new notation

$$\text{neighbours}(j) = \{1 \leq i \leq M : \boldsymbol{\theta}_i \text{ is a neighbour of } f_j\} \quad (6.7)$$

and (6.5) is re-expressible as  $p(\boldsymbol{\theta}, \mathbf{y}) = \prod_{j=1}^N f_j(\boldsymbol{\theta}_{\text{neighbours}(j)})$ .



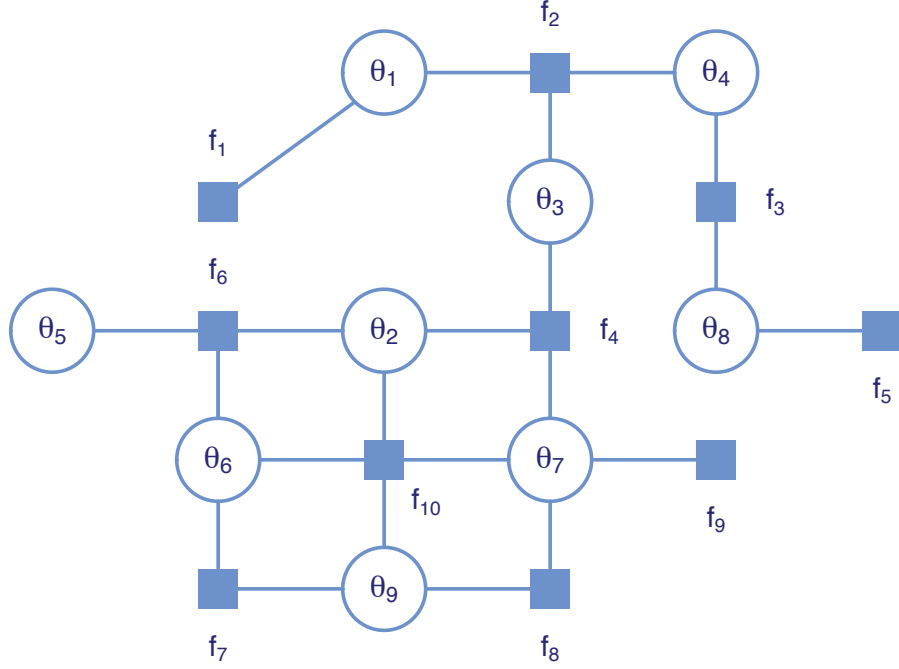


Figure 6.2: A factor graph corresponding to a Bayesian model with stochastic nodes  $\theta_1, \dots, \theta_9$  and factors  $f_1, \dots, f_{10}$ .

Wand (2015) provides a unified framework for obtaining the optimal  $q$ -densities through an iterative scheme concerned with updating two types of messages in a factor graph. Here we reiterate his scheme. For each  $1 \leq i \leq M$  and  $1 \leq j \leq N$ , the VMP updates for the first type of messages, *stochastic node to factor messages*, are

$$m_{\theta_i \rightarrow f_j}(\theta_i) \leftarrow \prod_{j' \neq j, i \in \text{neighbours}(j')} m_{f_{j'} \rightarrow \theta_i}(\theta_i) \quad (6.8)$$

and the second type of messages, *factor to stochastic node messages*, are

$$m_{f_j \rightarrow \theta_i}(\theta_i) \leftarrow \exp \left[ E_{f_j \rightarrow \theta_i} \left\{ \log f_j(\boldsymbol{\theta}_{\text{neighbours}(j)}) \right\} \right], \quad (6.9)$$

where the subscript  $m$  designates the nodes involved in the message passing and the direction in which the message is passed, and the notation  $E_{f_j \rightarrow \theta_i}$  denotes expectation with respect to the density function

$$\frac{1}{Z} \left( \prod_{i' \in \text{neighbours}(j) \setminus \{i\}} m_{f_j \rightarrow \theta_{i'}}(\theta_{i'}) m_{f_{i'} \rightarrow \theta_j}(\theta_{i'}) \right) \quad (6.10)$$

given  $Z$  is the corresponding normalising constant. The notation  $A \setminus B$  denotes the set

containing elements of  $A$  that are not in  $B$  if  $A$  and  $B$  are sets such that  $B \subseteq A$ . If  $\text{neighbours}(j) \setminus \{i\} = \emptyset$ , then the right hand side of 6.9 simplifies to  $f_j(\boldsymbol{\theta}_{\text{neighbours}(j)})$ . The convergence of these message updates is assessed by monitoring successive values of the lower bound of the marginal log-likelihood

$$\log \underline{p}(q; \mathbf{y}) = \sum_{i=1}^M E_{q(\boldsymbol{\theta}_i)} \{\log q(\boldsymbol{\theta}_i)\} + \sum_{j=1}^N E_q \{\log(f_j)\}. \quad (6.11)$$

The first term within the summation is known as the *entropy* or *differential entropy*. Wand (2015) provides a list of entropy expressions for the common exponential density functions. Upon convergence of the message updates (when the successive differences of (6.11) are very small), the optimal  $q$ -densities are obtained via

$$q^*(\boldsymbol{\theta}_i) \propto \prod_{i \in \text{neighbours}(j)} m_{f_j \rightarrow \boldsymbol{\theta}_i}(\boldsymbol{\theta}_i). \quad (6.12)$$

## 6.6 Illustrative example of simple mixed effects regression

Consider the random intercepts and slopes model (2.2),

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), & (6.13) \\ \mathbf{u} | \boldsymbol{\Sigma}^R &\sim N(\mathbf{0}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}^R), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), \\ \boldsymbol{\Sigma}^R | a_1^R, \dots, a_{q^R}^R &\sim \text{Inverse-Wishart} \left( \nu + q^R - 1, 2\nu \text{diag}(1/a_1^R, \dots, 1/a_{q^R}^R) \right), \\ a_r^R &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_{Rr}^2 \right), \quad r = 1, \dots, q^R \\ \sigma_\varepsilon^2 | a_\varepsilon &\sim \text{Inverse-Gamma} \left( \frac{1}{2}, 1/a_\varepsilon \right), \quad a_\varepsilon \sim \text{Inverse-Gamma} \left( \frac{1}{2}, 1/A_\varepsilon^2 \right), \end{aligned}$$

where the parameters are previously defined in Subsection 2.2.3. The joint posterior density function of model parameters and auxiliary variables is analytically intractable and numerically challenging, requiring variational approximations for fitting in practice. Mean field variational Bayes is a general class of variational methods for approximation of posterior density functions in the directed acyclic graph model infrastructure. Variational message passing is a message-passing version of the mean-field method, which leads to the identical posterior density function approximations since they are each founded upon the same Kullback-Leibler optimisation problem. Here we aim to obtain a mean-field approximation to the joint posterior density function:

$$p(\boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}^R, \sigma_\varepsilon^2, \mathbf{a}^R, a_\varepsilon) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\boldsymbol{\Sigma}^R) q(\sigma_\varepsilon^2) q(\mathbf{a}^R) q(a_\varepsilon).$$

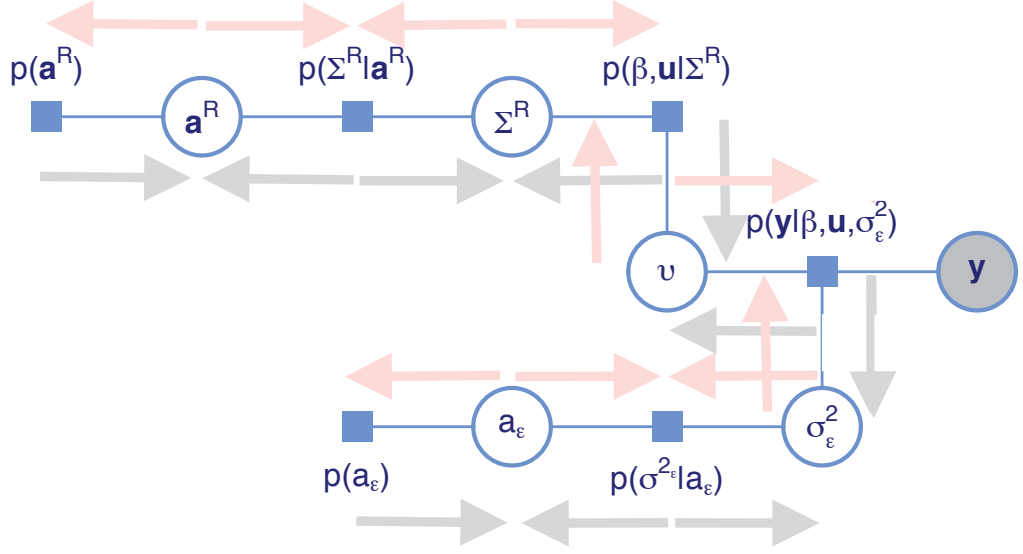


Figure 6.3: Each of the messages between neighbouring nodes on the factor graph for the random intercepts and slopes model illustrative example. The pink arrows depict the direction from the stochastic node to factor, while the grey arrows depict the direction from the factor to stochastic node.

VMP iterations for fitting model (6.13) involve updating 10 messages passed from factor to stochastic node (pink arrows) and 10 messages passed from stochastic node to factor (grey arrows) on the factor graph in Figure 6.3. Each message is a function of the stochastic node that is either received or sent by the message. Let  $\mathbf{v} \equiv [\boldsymbol{\beta}^\top \mathbf{u}^\top]^\top$  and  $\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$ , we first initialise each factor to stochastic node message to be an arbitrary density function:

$$\begin{aligned}
 m_{p(a_r^R) \rightarrow a_r^R}(a_r^R) &\propto \exp \left\{ \sum_{r=1}^{q^R} \begin{bmatrix} \log(a_r^R) \\ 1/a_r^R \end{bmatrix}^\top \boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} \right\}, \\
 m_{p(\Sigma^R | a_r^R) \rightarrow a_r^R}(a_r^R) &\propto \exp \left\{ \sum_{r=1}^{q^R} \begin{bmatrix} \log(a_r^R) \\ 1/a_r^R \end{bmatrix}^\top \boldsymbol{\eta}_{p(\Sigma^R | a_r^R) \rightarrow a_r^R} \right\}, \\
 m_{p(\Sigma^R | a^R) \rightarrow a^R}(\Sigma^R) &\propto \exp \left\{ \begin{bmatrix} \log |\Sigma^R| \\ \text{vec}\{(\Sigma^R)^{-1}\} \end{bmatrix}^\top \boldsymbol{\eta}_{p(\Sigma^R | a^R) \rightarrow \Sigma^R} \right\}, \\
 m_{p(\boldsymbol{\beta}, \mathbf{u} | \Sigma^R) \rightarrow \Sigma^R}(\Sigma^R) &\propto \exp \left\{ \begin{bmatrix} \log |\Sigma^R| \\ \text{vec}\{(\Sigma^R)^{-1}\} \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \Sigma^R) \rightarrow \Sigma^R} \right\}, \\
 m_{p(\boldsymbol{\beta}, \mathbf{u} | \Sigma^R) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}(\boldsymbol{\beta}, \mathbf{u}) &\propto \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \Sigma^R) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \right\},
 \end{aligned}$$

$$\begin{aligned}
 m_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}(\boldsymbol{\beta}, \mathbf{u}) &\propto \exp \left\{ \left[ \begin{array}{c} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^\top) \end{array} \right]^\top \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \right\}, \\
 m_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2) &\propto \exp \left\{ \left[ \begin{array}{c} \log(\sigma_\varepsilon^2) \\ 1/\sigma_\varepsilon^2 \end{array} \right]^\top \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \right\}, \\
 m_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}(\sigma_\varepsilon^2) &\propto \exp \left\{ \left[ \begin{array}{c} \log(\sigma_\varepsilon^2) \\ 1/\sigma_\varepsilon^2 \end{array} \right]^\top \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right\}, \\
 m_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon) &\propto \exp \left\{ \left[ \begin{array}{c} \log(a_\varepsilon) \\ 1/a_\varepsilon \end{array} \right]^\top \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} \right\}, \\
 m_{p(a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon) &\propto \exp \left\{ \left[ \begin{array}{c} \log(a_\varepsilon) \\ 1/a_\varepsilon \end{array} \right]^\top \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \right\},
 \end{aligned}$$

where the subscripted  $\boldsymbol{\eta}$  vectors are natural parameters for a particular exponential family density function and are initialised as, for example,

$$\boldsymbol{\eta}_{p(a_r^{\text{R}}) \rightarrow a_r^{\text{R}}} \leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_{\text{R}r}^2} \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}^{\text{R}}|a^{\text{R}}) \rightarrow a^{\text{R}}} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}.$$

Continuing this process leads to the first two and last four messages to be proportional to the inverse-Gamma distributions, the third and fourth messages are proportional to the inverse-Wishart distributions, and the fifth and sixth messages are proportional to the  $(p + mq^{\text{R}})$ -dimensional multivariate normal distributions.

The stochastic node to factor messages have the same functional forms as their reverse messages. The natural statistic vectors of messages remain unchanged and hence the VMP iterations boil down to computing and updating the corresponding natural parametric vectors. For example,

$$m_{a_r^{\text{R}} \rightarrow p(a_r^{\text{R}})}(a_r^{\text{R}}) \propto \exp \left\{ \sum_{r=1}^{q^{\text{R}}} \left[ \begin{array}{c} \log(a_r^{\text{R}}) \\ 1/a_r^{\text{R}} \end{array} \right]^\top \boldsymbol{\eta}_{a_r^{\text{R}} \rightarrow p(a_r^{\text{R}})} \right\}.$$

According to (6.8), the next step involves the updates of stochastic node to factor natural parameter vectors which have simple forms as follows:

$$\begin{aligned}
 \boldsymbol{\eta}_{a_r^{\text{R}} \rightarrow p(a_r^{\text{R}})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}^{\text{R}}|a_r^{\text{R}}) \rightarrow a_r^{\text{R}}}, & \boldsymbol{\eta}_{a_r^{\text{R}} \rightarrow p(\boldsymbol{\Sigma}^{\text{R}}|a_r^{\text{R}})} &\leftarrow \boldsymbol{\eta}_{p(a_r^{\text{R}}) \rightarrow a_r^{\text{R}}} \\
 \boldsymbol{\eta}_{\boldsymbol{\Sigma}^{\text{R}} \rightarrow p(\boldsymbol{\Sigma}^{\text{R}}|a^{\text{R}})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\text{R}}) \rightarrow \boldsymbol{\Sigma}^{\text{R}}}, & \boldsymbol{\eta}_{\boldsymbol{\Sigma}^{\text{R}} \rightarrow p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\text{R}})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}^{\text{R}}|a^{\text{R}}) \rightarrow \boldsymbol{\Sigma}^{\text{R}}} \\
 \boldsymbol{\eta}_{(\boldsymbol{\beta}, \mathbf{u}) \rightarrow p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\text{R}})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}, & \boldsymbol{\eta}_{(\boldsymbol{\beta}, \mathbf{u}) \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\text{R}}) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}
 \end{aligned}$$

6.6. ILLUSTRATIVE EXAMPLE OF SIMPLE MIXED EFFECTS REGRESSION

---

$$\begin{aligned} \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}, \quad \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}, \\ \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(a_\varepsilon)} \leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}, \quad \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon}. \end{aligned}$$

For model (6.13) each stochastic node to factor message is solely dependent upon one corresponding factor to stochastic node message due to the linear association between nodes as depicted in Figure 6.3. Based on (6.8) and (6.9),

$$m_{p(a_r^R) \rightarrow a_r^R} \propto p(a_r^R) \quad \text{and} \quad m_{p(a_\varepsilon) \rightarrow a_\varepsilon} \propto p(a_\varepsilon),$$

so the natural parameter updates for these two messages are simply

$$\boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} \leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_{Rr}^2} \end{bmatrix} \quad \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_\varepsilon^2} \end{bmatrix}$$

and remain constant throughout the iterations. The following step according to (6.9) is concerned with updating the factor to stochastic node messages. We focus on the stochastic node  $(\boldsymbol{\beta}, \mathbf{u})$  that involves updating

$$m_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}(\boldsymbol{\beta}, \mathbf{u}) \quad \text{and} \quad m_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^R) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}(\boldsymbol{\beta}, \mathbf{u}).$$

We begin with the first message

$$\begin{aligned} & m_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \tag{6.14} \\ & \propto \exp \left\{ \int_0^\infty m_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} m_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) d\sigma_\varepsilon^2 \right\} \\ & \propto \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^\top) \end{bmatrix}^\top \left[ \begin{array}{c} \mathbf{C}^\top \left( \int_0^\infty \frac{1}{\sigma_\varepsilon^2} p_{\text{IGnat}}(\sigma_\varepsilon^2, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \leftrightarrow \sigma_\varepsilon^2}) d\sigma_\varepsilon^2 \right) \mathbf{y} \\ -\frac{1}{2} \text{vec} \left\{ \mathbf{C}^\top \left( \int_0^\infty \frac{1}{\sigma_\varepsilon^2} p_{\text{IGnat}}(\sigma_\varepsilon^2, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \leftrightarrow \sigma_\varepsilon^2}) d\sigma_\varepsilon^2 \right) \mathbf{C} \right\} \end{array} \right] \right\} \end{aligned}$$

and then the second message

$$\begin{aligned} & m_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^R) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \tag{6.15} \\ & \propto \exp \left\{ \int_{\mathcal{R}_+^d} m_{\boldsymbol{\Sigma}^R \rightarrow p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^R)} m_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^R) \rightarrow \boldsymbol{\Sigma}^R} \log p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^R) d\boldsymbol{\Sigma}^R \right\} \end{aligned}$$

$$\propto \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^\top) \end{bmatrix}^\top \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left\{ \text{blockdiag} \left( \sigma_\beta^{-2} \mathbf{I}_p, \mathbf{I}_m \otimes \int_{\mathcal{R}_+^d} \frac{1}{\boldsymbol{\Sigma}^{\mathbf{R}}} \right) \right\} \\ \times p_{\text{IW}_{\text{nat}}}(\boldsymbol{\Sigma}^{\mathbf{R}}; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \leftrightarrow (\boldsymbol{\beta}, \mathbf{u})) d\boldsymbol{\Sigma}^{\mathbf{R}} \end{bmatrix} \right\}.$$

The expectation in (6.14) and (6.15) reduces to a linear combination of expected natural statistics. Using the algorithmic primitives defined in Section 6.3 then leads to

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \rightarrow (\boldsymbol{\beta}, \mathbf{u}) \leftarrow \begin{bmatrix} \mathbf{C}^\top \mathcal{I}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \leftrightarrow \sigma_\varepsilon^2) \mathbf{y} \\ -\frac{1}{2} \text{vec} \left( \mathbf{C}^\top \mathcal{I}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \leftrightarrow \sigma_\varepsilon^2) \mathbf{C} \right) \end{bmatrix}$$

and

$$\boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\mathbf{R}})} \rightarrow (\boldsymbol{\beta}, \mathbf{u}) \leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left\{ \text{blockdiag} \left( \sigma_\beta^{-2} \mathbf{I}_p, \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \leftrightarrow (\boldsymbol{\beta}, \mathbf{u})) \right) \right\} \end{bmatrix},$$

where the functions  $\mathcal{G}$  and  $\mathcal{I}$  are defined in Section 6.3. We repeat the same procedure for the remaining messages in the model. Once convergence of the messages has been attained, the  $q$ -density natural parameters can be obtained as:

$$\begin{aligned} \boldsymbol{\eta}_{q(\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\mathbf{R}})} \rightarrow (\boldsymbol{\beta}, \mathbf{u}) + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \rightarrow (\boldsymbol{\beta}, \mathbf{u}), \\ \boldsymbol{\eta}_{q(\sigma_\varepsilon^2)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2, \\ \boldsymbol{\eta}_{q(a_\varepsilon)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow a_\varepsilon + \boldsymbol{\eta}_{p(a_\varepsilon)} \rightarrow a_\varepsilon, \\ \boldsymbol{\eta}_{q(\boldsymbol{\Sigma}^{\mathbf{R}})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}^{\mathbf{R}}|a_r)} \rightarrow \boldsymbol{\Sigma}^{\mathbf{R}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}^{\mathbf{R}})} \rightarrow \boldsymbol{\Sigma}^{\mathbf{R}}, \\ \text{and } \boldsymbol{\eta}_{q(a_r^{\mathbf{R}})} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}^{\mathbf{R}}|a_r^{\mathbf{R}})} \rightarrow a_r^{\mathbf{R}} + \boldsymbol{\eta}_{p(a_r^{\mathbf{R}})} \rightarrow a_r^{\mathbf{R}}, \quad 1 \leq r \leq q^{\mathbf{R}}. \end{aligned}$$

Hence the resultant optimal  $q$ -densities in natural canonical form are:

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}) &\text{ is the } N_{\text{nat}}((\boldsymbol{\eta}_{q(\boldsymbol{\beta}, \mathbf{u})})_1, (\boldsymbol{\eta}_{q(\boldsymbol{\beta}, \mathbf{u})})_2, \text{vec}) \text{ density function,} \\ q^*(\sigma_\varepsilon^2) &\text{ is the Inverse-Gamma}_{\text{nat}}((\boldsymbol{\eta}_{q(\sigma_\varepsilon^2)})_1, (\boldsymbol{\eta}_{q(\sigma_\varepsilon^2)})_2) \text{ density function,} \\ q^*(a_\varepsilon) &\text{ is the Inverse-Gamma}_{\text{nat}}((\boldsymbol{\eta}_{q(a_\varepsilon)})_1, (\boldsymbol{\eta}_{q(a_\varepsilon)})_2) \text{ density function,} \\ q^*(\boldsymbol{\Sigma}^{\mathbf{R}}) &\text{ is the Inverse-Wishart}_{\text{nat}}((\boldsymbol{\eta}_{q(\boldsymbol{\Sigma}^{\mathbf{R}})})_1, (\boldsymbol{\eta}_{q(\boldsymbol{\Sigma}^{\mathbf{R}})})_2) \text{ density function,} \\ q^*(a_r^{\mathbf{R}}) &\text{ is the Inverse-Gamma}_{\text{nat}}((\boldsymbol{\eta}_{p(a_r^{\mathbf{R}})} \rightarrow a_r^{\mathbf{R}})_1, (\boldsymbol{\eta}_{p(a_r^{\mathbf{R}})} \rightarrow a_r^{\mathbf{R}})_2) \text{ density function.} \end{aligned}$$

Algorithm 17 gives the VMP iterative updates for obtaining optimal parameters for model

(6.13).

## 6.7 Arbitrarily large models viewpoint

VMP benefits from concepts such as natural parameter vectors and factor graphs, which makes it more amenable to modularisation than MFVB. The VMP algorithms can be easily extended to arbitrarily large models via the notion of *factor graph fragments* that were first introduced in Wand (2015):

**Definition 6.7.1.** A factor graph fragment is a sub-graph of a factor graph consisting of a single factor and each of its neighbouring stochastic nodes.

Wand (2015) described the factor graph fragments as building blocks for a wide range of VMP-based Bayesian semiparametric mixed models involving, for example, multiple predictors and interactions, penalised splines or wavelets, higher-level random effects and non-Gaussian responses. Interested readers would benefit from reading Section 4 of his paper on derivation of the five fundamental fragments. With the help of these fragment definitions, the updates of natural parameters for the factor to stochastic node messages only need to be derived and coded once and then be integrated into different compartmentalised functions. This potentially saves a huge amount of time in terms of derivation and programming.

Consider an example where we extend the simple mixed model (6.13) to a three-level semiparametric mixed model with Gaussian response

$$\begin{aligned}
 \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}), \quad \boldsymbol{\beta} \sim N(0, \sigma_\beta^2 \mathbf{I}_p), \\
 \mathbf{u} | \sigma_{u\ell}^2, \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}} &\sim N\left(\mathbf{0}, \begin{bmatrix} \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I}_{q_\ell^G}) & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left[ \begin{array}{cc} \boldsymbol{\Sigma}^{\text{RL3}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_i} \otimes \boldsymbol{\Sigma}^{\text{RL2}} \end{array} \right] \end{bmatrix}\right), \\
 \boldsymbol{\Sigma}^{\text{RL2}} | a_1^{\text{RL2}}, \dots, a_{q^{\text{RL2}}}^{\text{RL2}} &\sim \text{Inverse-Wishart}\left(\nu + q^{\text{RL2}} - 1, 2\nu \text{diag}(1/a_1^{\text{RL2}}, \dots, 1/a_{q^{\text{RL2}}}^{\text{RL2}})\right), \\
 \boldsymbol{\Sigma}^{\text{RL3}} | a_1^{\text{RL3}}, \dots, a_{q^{\text{RL3}}}^{\text{RL3}} &\sim \text{Inverse-Wishart}\left(\nu + q^{\text{RL3}} - 1, 2\nu \text{diag}(1/a_1^{\text{RL3}}, \dots, 1/a_{q^{\text{RL3}}}^{\text{RL3}})\right), \\
 a_{r_{L2}}^{\text{RL2}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{RL2}}^2\right), \quad r_{L2} = 1, \dots, q^{\text{RL2}}, \\
 a_{r_{L3}}^{\text{RL3}} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{RL3}}^2\right), \quad r_{L3} = 1, \dots, q^{\text{RL3}}, \\
 \sigma_{u\ell}^2 | a_{u\ell} &\stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{u\ell}\right), \quad a_{u\ell} \stackrel{\text{ind.}}{\sim} \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{u\ell}^2\right),
 \end{aligned} \tag{6.16}$$

where the model parameters are previously defined in Subsection 3.7.1. The joint posterior density function of all random variables in model (6.16) admits the factorisation

$$p(\boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \sigma_u^2, \sigma_\varepsilon^2, \mathbf{a}^{\text{RL2}}, \mathbf{a}^{\text{RL3}}, \mathbf{a}_u, a_\varepsilon) = \tag{6.17}$$

**Initialise all factor to stochastic node messages to be arbitrary density functions:**

$$\begin{aligned}
 \boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} &\leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_{Rr}^2} \end{bmatrix} ; \boldsymbol{\eta}_{p(\Sigma^R|a_r^R) \rightarrow a_r^R} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} ; \boldsymbol{\eta}_{p(\Sigma^R|a^R) \rightarrow \Sigma^R} \leftarrow \begin{bmatrix} -3 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_2) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow \Sigma^R} &\leftarrow \begin{bmatrix} -3 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_2) \end{bmatrix} ; \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow (\beta, \mathbf{u})} \leftarrow \begin{bmatrix} -\mathbf{0}_{p+mq^R} \\ -\frac{1}{2}\text{vec}(\mathbf{I}_{p+mq^R}) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\beta, \mathbf{u})} &\leftarrow \begin{bmatrix} -\mathbf{0}_{p+mq^R} \\ -\frac{1}{2}\text{vec}(\mathbf{I}_{p+mq^R}) \end{bmatrix} ; \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} ; \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_\varepsilon^2} \end{bmatrix}
 \end{aligned}$$

**Cycle through updates:**

(a) For stochastic node to factor messages:

For  $r = 1, \dots, q^R$

$$\begin{aligned}
 \boldsymbol{\eta}_{a_r^R \rightarrow p(a_r^R)} &\leftarrow \boldsymbol{\eta}_{p(\Sigma^R|a_r^R) \rightarrow a_r^R} ; \boldsymbol{\eta}_{a_r^R \rightarrow p(\Sigma^R|a_r^R)} \leftarrow \boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} \\
 \boldsymbol{\eta}_{\Sigma^R \rightarrow p(\Sigma^R|a^R)} &\leftarrow \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow \Sigma^R} ; \boldsymbol{\eta}_{\Sigma^R \rightarrow p(\beta, \mathbf{u}|\Sigma^R)} \leftarrow \boldsymbol{\eta}_{p(\Sigma^R|a^R) \rightarrow \Sigma^R} \\
 \boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow p(\beta, \mathbf{u}|\Sigma^R)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\beta, \mathbf{u})} ; \boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2)} \leftarrow \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow (\beta, \mathbf{u})} \\
 \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} ; \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \\
 \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(a_\varepsilon)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} ; \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon}
 \end{aligned}$$

(b) For factor to stochastic node messages:

For  $r = 1, \dots, q^R$

$$\begin{aligned}
 \boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} &\leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_{Rr}^2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\Sigma^R|a_r^R) \rightarrow a_r^R} &\leftarrow \begin{bmatrix} -\frac{1}{2}(\nu + q^R - 1) \\ -\nu \mathcal{G} \left( \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow \Sigma^R} + \boldsymbol{\eta}_{p(\Sigma^R|a^R) \rightarrow \Sigma^R} \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\Sigma^R|a^R) \rightarrow \Sigma^R} &\leftarrow \begin{bmatrix} -\frac{1}{2}(A_{\Sigma^R} + q^R - 1) \\ -\frac{1}{2} \sum_{r=1}^2 \text{vec} \left\{ 2\nu \text{diag} \left( \mathcal{I} \left( \boldsymbol{\eta}_{p(\Sigma^R|a_r^R) \rightarrow a_r^R} + \boldsymbol{\eta}_{p(a_r^R) \rightarrow a_r^R} \right) \right) \right\} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma^R) \rightarrow \Sigma^R} &\leftarrow \begin{bmatrix} -\frac{1}{2} \sum_{i=1}^m n_i \\ \mathcal{F} \left( \boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2)} + \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\beta, \mathbf{u})} \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow (\beta, \mathbf{u})} &\leftarrow \begin{bmatrix} \mathbf{C}^\top \mathcal{I} \left( \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \right) \mathbf{y} \\ \frac{1}{2} \text{vec} \left\{ \mathbf{C}^\top \mathcal{I} \left( \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \right) \mathbf{C} \right\} \end{bmatrix}
 \end{aligned}$$



$$\begin{aligned}
 \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\Sigma^{\mathbf{R}}) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \begin{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} \\ -\frac{1}{2} \text{vec} \left( \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\Sigma^{\mathbf{R}}|a^{\mathbf{R}}) \rightarrow \Sigma^{\mathbf{R}}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\Sigma^{\mathbf{R}}) \rightarrow \Sigma^{\mathbf{R}}}) \end{bmatrix} \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) \rightarrow \sigma_{\varepsilon}^2} &\leftarrow \begin{bmatrix} -\frac{1}{2} \sum_{i=1}^m n_i \\ -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^{n_i} \text{vec} \left( \mathcal{H}(\boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\Sigma^{\mathbf{R}}) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}; \mathbf{y}_{ij}; \mathbf{C}_{ij}) \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow \sigma_{\varepsilon}^2} &\leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\mathcal{I} \left( \boldsymbol{\eta}_{p(a_{\varepsilon}) \rightarrow a_{\varepsilon}} + \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow a_{\varepsilon}} \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(a_{\varepsilon}) \rightarrow a_{\varepsilon}} &\leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{1}{A_{\varepsilon}^2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow a_{\varepsilon}} &\leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\mathcal{I} \left( \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) \rightarrow \sigma_{\varepsilon}^2} + \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow \sigma_{\varepsilon}^2} \right) \end{bmatrix}
 \end{aligned}$$

until the increase in  $\log p(\mathbf{y}; q)$  is negligible.

**Obtain  $q$ -density natural parameters via summation of messages:**

$$\boldsymbol{\eta}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\Sigma^{\mathbf{R}}) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}$$

$$\boldsymbol{\eta}_{q(\sigma_{\varepsilon}^2)} \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) \rightarrow \sigma_{\varepsilon}^2} + \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow \sigma_{\varepsilon}^2}$$

$$\boldsymbol{\eta}_{q(a_{\varepsilon})} \leftarrow \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon}) \rightarrow a_{\varepsilon}} + \boldsymbol{\eta}_{p(a_{\varepsilon}) \rightarrow a_{\varepsilon}}$$

$$\boldsymbol{\eta}_{q(\Sigma^{\mathbf{R}})} \leftarrow \boldsymbol{\eta}_{p(\Sigma^{\mathbf{R}}|a^{\mathbf{R}}) \rightarrow \Sigma^{\mathbf{R}}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u}|\Sigma^{\mathbf{R}}) \rightarrow \Sigma^{\mathbf{R}}}$$

For  $r = 1, \dots, q^{\mathbf{R}}$

$$\boldsymbol{\eta}_{q(a_r^{\mathbf{R}})} \leftarrow \boldsymbol{\eta}_{p(\Sigma^{\mathbf{R}}|a_r^{\mathbf{R}}) \rightarrow a_r^{\mathbf{R}}} + \boldsymbol{\eta}_{p(a_r^{\mathbf{R}}) \rightarrow a_r^{\mathbf{R}}}$$

---

Algorithm 17: Variational message passing iterative scheme for obtaining the optimal parameters for model (6.13).

## 6.7. ARBITRARILY LARGE MODELS VIEWPOINT

$$\begin{aligned}
& p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) p(\boldsymbol{\beta}, \mathbf{u} | (\boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \sigma_u^2)) p(\boldsymbol{\Sigma}^{\text{RL2}} | \mathbf{a}^{\text{RL2}}) p(\boldsymbol{\Sigma}^{\text{RL3}} | \mathbf{a}^{\text{RL3}}) \\
& \times p(\sigma_u^2 | \mathbf{a}_u) p(\sigma_\varepsilon^2 | a_\varepsilon) p(\mathbf{a}^{\text{RL2}}) p(\mathbf{a}^{\text{RL3}}) p(\mathbf{a}_u) p(a_\varepsilon).
\end{aligned}$$

Treating (6.17) as a set of parameter functions corresponding to each factor gives the factor graph in Figure 6.4. This factor graph has 10 graph fragments, with five of them have the same form as the fragments of Figure 6.3 and are coloured in light blue. The darker blue fragments are the additional ones.

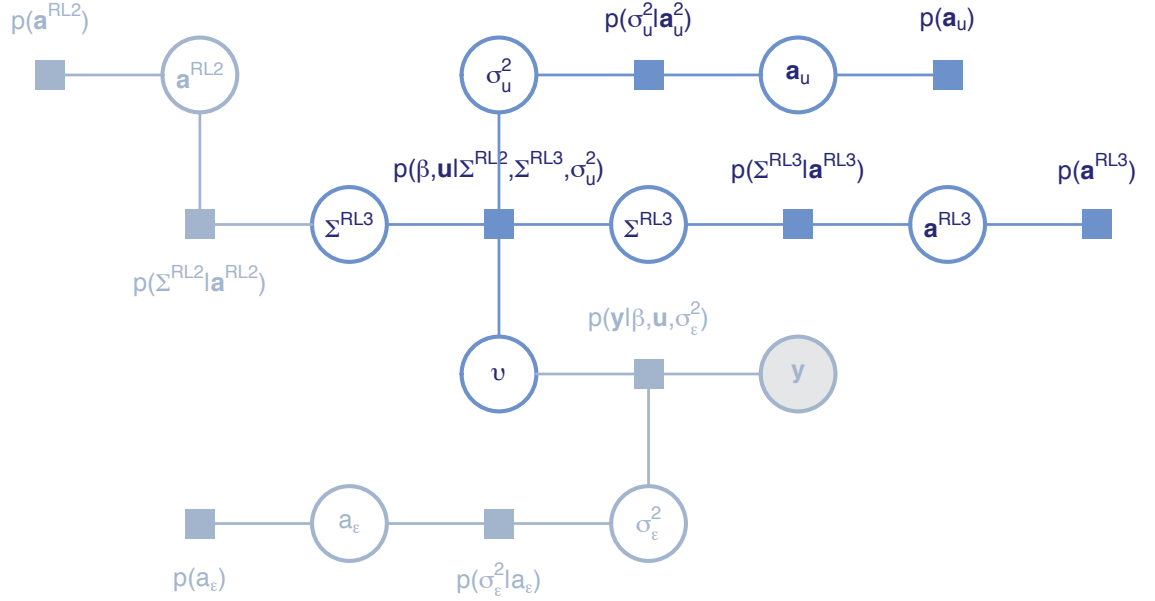


Figure 6.4: Diagrammatic depiction of the extension from simple mixed effects regression to three-level penalised spline mixed regression. The fragments shown in darker blue are the ones that are additional to the fragments in the simple mixed model. The fragments shown in grey are presented in Figure 6.3.

The stochastic node to factor messages in Figure 6.4 have trivial updates analogous to those given in Section 6.6. The factor to stochastic node messages are however more complicated, but the majority of the darker blue fragments have identical or very similar updates as the ones given in factor graph Figure 6.3. For example, the messages passed from  $p(\sigma_u^2 | \mathbf{a}_u)$  to  $\sigma_u^2$  and  $\mathbf{a}_u$  have the same form as those passed from  $p(\sigma_\varepsilon^2)$  to  $\sigma_\varepsilon^2$  and  $a_\varepsilon$  for model (6.13). The natural parameter updates  $\boldsymbol{\eta}_{p(\sigma_u^2 | \mathbf{a}_u)} \rightarrow \sigma_u^2$  and  $\boldsymbol{\eta}_{p(\sigma_u^2 | \mathbf{a}_u)} \rightarrow \mathbf{a}_u$  take the same forms as those for  $\boldsymbol{\eta}_{p(\sigma_\varepsilon^2 | a_\varepsilon)} \rightarrow \sigma_\varepsilon^2$  and  $\boldsymbol{\eta}_{p(\sigma_\varepsilon^2 | a_\varepsilon)} \rightarrow a_\varepsilon$  given in Section 6.6. Similarly, the message passed from  $p(\mathbf{a}_u)$  to  $\mathbf{a}_u$  has the same form as that passed from  $p(a_\varepsilon)$  to  $a_\varepsilon$ . The natural parameter updates  $\boldsymbol{\eta}_{p(\mathbf{a}_u)} \rightarrow \mathbf{a}_u$  takes the same form as that for  $\boldsymbol{\eta}_{p(a_\varepsilon)} \rightarrow a_\varepsilon$ . Similar arguments can be applied to the messages passed from  $p(\boldsymbol{\Sigma}^{\text{RL3}} | \mathbf{a}^{\text{RL3}})$  to  $\boldsymbol{\Sigma}^{\text{RL3}}$  and  $\mathbf{a}^{\text{RL3}}$ , and messages from  $p(\mathbf{a}^{\text{RL3}})$  to  $\mathbf{a}^{\text{RL3}}$ . It remains to take care of the darker blue fragments that are at the centre of the factor graph. The message passed from

## 6.8. CONCLUDING REMARKS

---

$p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2)$  to  $\boldsymbol{\sigma}_u^2$  is

$$\begin{aligned} m_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow \sigma_{u\ell}^2}(\sigma_{u\ell}^2) \\ = \exp \left\{ \sum_{\ell=1}^L \begin{bmatrix} \log(\sigma_{u\ell}^2) \\ 1/\sigma_{u\ell}^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow \sigma_{u\ell}^2} \right\}, \end{aligned}$$

for  $\ell = 1, \dots, L$ , where

$$\boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow \sigma_{u\ell}^2} \leftarrow \begin{bmatrix} -\frac{1}{2}q_\ell^{\text{G}} \\ \mathcal{F} \left( \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\sigma}_\varepsilon^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \right) \\ + \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \end{bmatrix}.$$

The message passed from  $p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2)$  to  $(\boldsymbol{\beta}, \mathbf{u})$  is

$$\begin{aligned} m_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})}(\boldsymbol{\beta}, \mathbf{u}) \\ = \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} \right\}, \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \rightarrow (\boldsymbol{\beta}, \mathbf{u})} &\leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left\{ \text{blockdiag} \left( \sigma_\beta^{-2} \mathbf{I}, \boldsymbol{\Omega}_\Sigma, \boldsymbol{\Omega}_u \right) \right\} \end{bmatrix}, \\ \boldsymbol{\Omega}_\Sigma &= \text{blockdiag}_{1 \leq i \leq m} \begin{bmatrix} \mathbf{I}_{n_i} \otimes \mathcal{F}(\boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \leftrightarrow \boldsymbol{\Sigma}^{\text{RL2}}}) & \mathbf{0} \\ \mathbf{0} & \mathcal{F}(\boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \leftrightarrow \boldsymbol{\Sigma}^{\text{RL3}}}) \end{bmatrix}, \\ \boldsymbol{\Omega}_u &= \text{blockdiag}_{1 \leq i \leq m} \left( \mathcal{I} \left( \boldsymbol{\eta}_{p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2) \leftrightarrow \sigma_{u\ell}^2} \right) \mathbf{I}_{q_\ell^{\text{G}}} \right), \end{aligned}$$

where the functions  $\mathcal{F}$  and  $\mathcal{I}$  are defined in Section 6.3. The message passed from  $p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{RL2}}, \boldsymbol{\Sigma}^{\text{RL3}}, \boldsymbol{\sigma}_u^2)$  to  $\boldsymbol{\Sigma}^{\text{RL2}}$  or  $\boldsymbol{\Sigma}^{\text{RL3}}$  takes the same form as that passed from  $p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}^{\text{R}})$  to  $\boldsymbol{\Sigma}^{\text{R}}$  in Section 6.6.

## 6.8 Concluding remarks

We demonstrated that the notion of message passing can be used to streamline algebra calculations and programming for fast approximate inference in large Bayesian Gaussian semiparametric mixed models. The resultant VMP algorithms have ready-to-implement closed-form expressions and are easier to extend to a broad class of arbitrarily large

## 6.8. CONCLUDING REMARKS

---

Bayesian mixed models, via factor graph fragments, than the MFVB approach.

## Chapter 7

# Conclusions and Future Directions

*The goal is to transform data into information, information into insight.*

Carly Fiorina

We have exhaustively exemplified the use of variational methods in the setting of Bayesian longitudinal and multilevel data analysis. In this thesis we introduced sophisticated variational inference techniques for fitting and inference in a wide variety of Bayesian semiparametric mixed models. Four common response distributions, including Gaussian, Student- $t$ , Bernoulli and Poisson, were explored. MFVB is a principled approach to approximate inference, where it defines a family of approximating distributions whose members are simpler in some sense than the actual posterior and then seeks the optimal member of that family to minimising the Kullback-Leibler divergence (the dissimilarity criterion) between the approximating distribution and the exact posterior. Through a series of numerical studies and real data examples, we demonstrated that MFVB approximations have attractive computation and accuracy trade-offs. To further improve speed and memory efficiency of the naïve/direct MFVB algorithms we used matrix permutation and block decomposition to streamline the estimation of the large sparse effects covariance matrix. Rather than deciding whether this new methodology dominates the traditional ones in general, we discussed the strengths and weaknesses of the variational and standard MCMC methods. While MCMC sampling provides theoretical guarantees of accuracy, variational methods, in general, suffer some accuracy loss, albeit minimal, by the nature of the restriction on which they are based. In several instances, variational approximations can be shown to be overconfident, i.e. they underestimate posterior variances (e.g. Bishop, 2006; Wang and Titterington, 2006). This raises the question of whether this underestimation

---

affects model selection when variance components are a focus? As previously shown in You *et al.* (2014) such criticism does not hold for general variational methods in an asymptotic sense. The authors have proven the variational estimators for linear regression models achieve desirable frequentist properties such as consistency and can obtain asymptotically valid standard errors under mild regularity conditions. In addition, Ormerod *et al.* (2014) extended You *et al.* (2014)'s approach by inducing sparsity on the linear regression coefficients via a spike and slab prior and showed that the resultant estimates are consistent with valid standard errors, and that the true model can be found at an exponential rate in  $n$ . Extending these theoretical results to non-Gaussian responses and other data scenarios was not considered in the above papers but would be worthwhile to investigate in the future. Nonetheless, variational algorithms provide impressively fast and analytical approximations to otherwise unattainable posteriors. For example, we found that variational methods can be exploited to yield closed-form expressions that approximate the posteriors for the logistic regression effect parameters. It is good to keep in mind that both MCMC and variational methods are computational paradigms, each offers a range of different algorithmic approaches which trade off between speed, accuracy and ease of implementation.

The improvements in speed and memory efficiency offered by MFVB are not a simple novelty. Putting accuracy considerations aside, the algebraic infrastructure we laid out in this thesis has far-reaching implications beyond the examples considered here, for the analysis of big datasets via Bayesian hierarchical models as both continue to grow in size and complexity. We showed the versatility and utility of MFVB algorithms in terms of speed for moderately large data, and easy extendability to complicated, but more realistic scenarios, including semiparametric regression, higher-level mixed models, measurement error and missing data models, models with group-specific curves, and real-time/online processing for high velocity data. The contributions made in this thesis open up a broader class of models considering variational methods as a strong complement to the more traditional MCMC methods for exploratory data analysis and iterative model fitting. Furthermore, in order to get the best of both worlds, the MFVB approximation can always be used to initialise a sampling based approach which tends to asymptotically recover the actual posterior. Thus the fast variational methods can be used to quickly explore and respecify models which then lead to improved convergence of the asymptotically exact MCMC sampling algorithms.

Variational inference opens the door to several promising research directions. The MFVB algorithms derived in this thesis provide an efficient and convenient iterative scheme for updating the variational parameters, but in most examples we used conjugate priors, often chosen for reasons of tractability. If  $q_i(\boldsymbol{\theta}_i)$  does not belong to a recognisable density family, some optimisation techniques are required to estimate the marginal likeli-

---

hood. Future work can move beyond closed-form updates and fully factorised approximate posteriors, as well as allow arbitrary priors for the hyperparameters (e.g. Neville *et al.*, 2014). For examples, Teh *et al.* (2006) introduce *collapsed variational inference*, which marginalises out some of the hidden variables, training simple closed-form updates for low-dimensional posteriors. Hoffman *et al.* (2013) develop *stochastic variational inference* that facilitates stochastic optimisation to optimise the variational objective. The authors use structured variational distributions to relax the mean-field restriction for better approximating complex posteriors such as those arising in time-series models. Tan and Nott (2014) extend the stochastic variational inference approach of Hoffman *et al.* (2013) by combining non-conjugate variational messaging passing with algorithms from stochastic optimisation which work with mini-batches of data, with applications to non-conjugate generalised linear mixed models. Related work on non-conjugate models by Wang and Blei (2013) use the Laplace approximation for updating non-conjugate variational factors and, as described in the previous chapters, Knowles and Minka (2011) introduce the non-conjugate variational message passing framework for variational approximations in the exponential family. Finally, the popular divide and recombine strategy introduced in Nott *et al.* (2013) involves partitioning a large dataset into smaller subsets and then combining the approximating distributions that have been learned in parallel from each separate subset using the *hybrid variational Bayes algorithm*, a combination of the mean-field and fixed-form variational methods.

We live in an increasingly data-rich world, with ever greater desires and responsibilities to extract signals from noisy data. Even though the simulated and real data examples considered in this thesis are not gigabytes in size, coupled with modern data subsetting techniques, the variational inference framework presented in this thesis offers the opportunity to develop many fast algorithms that can be applied to even larger datasets which require storage greater than what is possible on a standard computer.

# References

- Abrevaya, J. (2006). Estimating the effect of smoking on birth outcomes using a matched panel data approach. *Journal of Applied Econometrics*, **21**:489–519.
- Albert, J. H. and Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American statistical Association*, **88**:669–679.
- Baicker, K., Buckles, K. S., and Chandra, A. (2006). Geographic variation in the appropriate use of cesarean delivery. *Health Affairs*, **25**:w355–w367.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- Bishop, C. M. (2008). A new framework for machine learning. Em *Computational Intelligence: Research Frontiers*, pgs. 1–24. Springer.
- Bishop, C. M., Spiegelhalter, D., and Winn, J. (2002). VIBES: a variational inference engine for Bayesian networks. Em *Advances in neural information processing systems*, pgs. 777–784.
- Bragg, F., Cromwell, D. A., Edozien, L. C., Gurol-Urganci, I., Mahmood, T. A., Templeton, A., and van der Meulen, J. H. (2010). Variation in rates of caesarean section among english NHS trusts after accounting for maternal and clinical risk: cross sectional study. *British Medical Journal*, **341**:c5065.
- Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, **88**:9–25.
- Brumback, B. A. and Rice, J. A. (1998). Smoothing spline models for the analysis of nested and crossed samples of curves. *Journal of the American Statistical Association*, **93**:961–976.
- Carroll, R. J., Ruppert, D., Stefanski, L. A., and Crainiceanu, C. M. (2006). *Measurement error in nonlinear models: a modern perspective*. CRC press.
- Challis, E. and Barber, D. (2013). Gaussian Kullback-Leibler approximate inference. *The*



## REFERENCES

---

- Journal of Machine Learning Research*, **14**:2239–2286.
- Chappell, M., Groves, A. R., Whitcher, B., and Woolrich, M. W. (2009). Variational Bayesian inference for a nonlinear forward model. *Signal Processing, IEEE Transactions on*, **57**:223–236.
- Chen, H. and Wang, Y. (2011). A penalized spline approach to functional mixed effects model analysis. *Biometrics*, **67**:861–870.
- Clayton, D. G. (1992). Models for the analysis of cohort and case-control studies with inaccurately measured exposures. *Statistical models for longitudinal studies of health*, pgs. 301–331.
- Coull, B. A., Catalano, P. J., and Godleski, J. J. (2000). Semiparametric analyses of cross-over data with repeated measures. *Journal of Agricultural, Biological, and Environmental Statistics*, **5**:417–429.
- Daniels, M. J. and Hogan, J. W. (2008). *Missing data in longitudinal studies: Strategies for Bayesian modeling and sensitivity analysis*. CRC Press.
- Dechter, R. and Pearl, J. (1988). *Network-based heuristics for constraint-satisfaction problems*. New York: Springer.
- Diggle, P. J., Heagerty, P., Liang, K.-Y., and Zeger, S. (2002). *Analysis of longitudinal data*. Oxford: Oxford University Press.
- Durbán, M., Harezlak, J., Wand, M. P., and Carroll, R. J. (2005). Simple fitting of subject-specific curves for longitudinal data. *Statistics in Medicine*, **24**:1153–1167.
- Faes, C., Ormerod, J. T., and Wand, M. P. (2011). Variational Bayesian inference for parametric and nonparametric regression with missing data. *Journal of the American Statistical Association*, **106**:959–971.
- Fitzmaurice, G. M., Davidian, M., Verbeke, G., and Molenberghs, G. (2008). *Longitudinal data analysis*. CRC Press.
- Fitzmaurice, G. M., Laird, N. M., and Ware, J. H. (2012). *Applied longitudinal analysis*, volume 998. New Jersey: John Wiley & Sons.
- Fong, Y., Rue, H., and Wakefield, J. (2009). Bayesian inference for generalized linear mixed models. *Biostatistics*, pg. kxp053.
- Frey, B. J. (1998). *Graphical models for machine learning and digital communication*. Massachusetts: MIT press.

## REFERENCES

---

- Gelfand, A. E. and Smith, A. F. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association*, **85**(410):398–409.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, **1**:515–533.
- Gelman, A. and Hill, J. (2007). *Data analysis using regression and multilevel or hierarchical models*. New York: Cambridge University Press.
- Geman, S. and Geman, D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, **(6)**:721–741.
- Goldstein, H. (1986). Multilevel mixed linear model analysis using iterative generalized least squares. *Biometrika*, **73**:43–56.
- Goldstein, H. (2011). *Multilevel statistical models*, volume 922. New Jersey: John Wiley & Sons.
- Harville, D. A. (2008). *Matrix algebra from a statistician's perspective*, volume 1. New York: Springer.
- Hastie, T., Tibshirani, R., Friedman, J., and Franklin, J. (2005). The elements of statistical learning: data mining, inference and prediction. *The Mathematical Intelligencer*, **27**:83–85.
- Hastie, T. J. and Tibshirani, R. J. (1990). *Generalized additive models*, volume 43. CRC Press.
- Hoffman, M., Bach, F. R., and Blei, D. M. (2010). Online learning for latent Dirichlet allocation. In *advances in neural information processing systems*, pgs. 856–864.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, **14**:1303–1347.
- Huang, A. and Wand, M. P. (2013). Simple marginally noninformative prior distributions for covariance matrices. *Bayesian Analysis*, **8**:439–452.
- Jaakkola, T. S. and Jordan, M. I. (1997). A variational approach to Bayesian logistic regression models and their extensions. In *Sixth International Workshop on Artificial Intelligence and Statistics*.
- Jaakkola, T. S. and Jordan, M. I. (2000). Bayesian parameter estimation via variational

## REFERENCES

---

- methods. *Statistics and Computing*, **10**:25–37.
- Jordan, M. I. (2004). Graphical models. *Statistical Science*, **19**:140–155.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1999). An introduction to variation methods for graphical models. *Machine Learning*, **37**:183–233.
- Knowles, D. A. and Minka, T. (2011). Non-conjugate variational message passing for multinomial and binary regression. Em *Advances in Neural Information Processing Systems*, pgs. 1701–1709.
- Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data. *Biometrics*, **38**:963–974.
- Lee, C. Y. Y., Homer, C., Bisits, A., and Ryan, L. (2015). Increasing variation in hospital caesarean section rates among low-risk nulliparous women in australia, from 1994 to 2010. *Birth in revision*.
- Lee, C. Y. Y. and Wand, M. P. (2015a). Streamlined mean field variational Bayes for longitudinal and multilevel data analysis. *Biometrical Journal in press*.
- Lee, C. Y. Y. and Wand, M. P. (2015b). Variational methods for fitting complex Bayesian mixed effects models to health data. *Statistics in Medicine*. DOI 10.1002/sim.6737.
- Li, Y. and Ruppert, D. (2008). On the asymptotics of penalized splines. *Biometrika*, **95**:415–436.
- Ligges, U., Thomas, A., Spiegelhalter, D., Best, N., Lunn, D., Rice, K., and Sturtz, S. (2009). *BRugs 0.5: OpenBUGS and its R/S-PLUS interface BRugs*.
- Little, R. J. and Rubin, D. B. (2014). *Statistical analysis with missing data*. New Jersey: John Wiley & Sons.
- Luts, J. (2015). Real-time semiparametric regression for distributed data sets. *Knowledge and Data Engineering, IEEE Transactions on*, **27**:545–557.
- Luts, J., Broderick, T., and Wand, M. P. (2014). Real-time semiparametric regression. *Journal of Computational and Graphical Statistics*, **23**:589–615.
- Magnus, J. R. and Neudecker, H. (1995). *Matrix differential calculus with applications in statistics and econometrics*. New Jersey: John Wiley & Sons.
- Maringwa, J. T., Geys, H., Shkedy, Z., Faes, C., Molenberghs, G., Aerts, M., Ammel, K. V., Teisman, A., and Bijmens, L. (2008). Application of semiparametric mixed models and simultaneous confidence bands in a cardiovascular safety experiment

## REFERENCES

---

- with longitudinal data. *Journal of Biopharmaceutical Statistics*, **18**:1043–1062.
- Marley, J. and Wand, M. P. (2010). Non-standard semiparametric regression via **brugs**. *Journal of Statistical Software*, **37**:1–30.
- McGrory, C. A. and Titterton, D. (2007). Variational approximations in Bayesian model selection for finite mixture distributions. *Computational Statistics and Data Analysis*, **51**:5352–5367.
- Menictas, M. (2015). *Variational inference for heteroscedastic and longitudinal regression models*. PhD thesis, University of Technology Sydney, School of Mathematical and Physical Sciences.
- Menictas, M. and Wand, M. P. (2013). Variational inference for marginal longitudinal semiparametric regression. *Stat*, **2**:61–71.
- Minka, T. (2005). Divergence measures and message passing. *Technical report, Microsoft Research*.
- Minka, T., Winn, J., Guiver, J., and Kannan, A. (2009). *Infer.NET 2.3*. Microsoft Research Cambridge, Cambridge, UK.
- Neville, S. E., Ormerod, J. T., and Wand, M. P. (2014). Mean field variational Bayes for continuous sparse signal shrinkage: pitfalls and remedies. *Electronic Journal of Statistics*, **8**:1113–1151.
- Nott, D. J., Tran, M.-N., Kuk, A. Y., and Kohn, R. (2013). Efficient variational inference for generalized linear mixed models with large datasets. *Unpublished*.
- Opper, M. and Archambeau, C. (2009). The variational Gaussian approximation revisited. *Neural Computation*, **21**:786–792.
- Organisation for Economic Co-Operation and Development (2011). *Caesarean sections. In: Health at a glance 2011: OECD Indicators*. OECD Publishing.
- Ormerod, J. T. and Wand, M. P. (2010). Explaining variational approximation. *The American Statistician*, **64**:140–153.
- Ormerod, J. T. and Wand, M. P. (2012). Gaussian variational approximate inference for generalized linear mixed models. *Journal of Computational and Graphical Statistics*, **21**:2–17.
- Ormerod, J. T., You, C., and Muller, S. (2014). A variational Bayes approach to variable selection. (Unpublished).

## REFERENCES

---

- O’Sullivan, F. (1986). A statistical perspective on ill-posed inverse problems. *Statistical Science*, **1**:502–527.
- Pham, T. H., Ormerod, J. T., and Wand, M. P. (2013). Mean field variational Bayesian inference for nonparametric regression with measurement error. *Computational Statistics and Data Analysis*, **68**:375–387.
- Pinheiro, J., Bates, D., DebRoy, S., and Sarkar, D. (2014). Mixed GAM computation vehicle with GCV/AIC/REML smoothness estimation. *R package version 1.8*.
- R Development Core Team (2015). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rabe-Hesketh, S. and Skrondal, A. (2008). *Multilevel and longitudinal modeling using Stata*. Texas: STATA press.
- Rao, C. R. (1972). Estimation of variance and covariance components in linear models. *Journal of the American Statistical Association*, **67**:112–115.
- Raudenbush, S. W., Yang, M.-L., and Yosef, M. (2000). Maximum likelihood for generalized linear models with nested random effects via high-order, multivariate laplace approximation. *Journal of computational and Graphical Statistics*, **9**:141–157.
- Rhode, D. and Wand, M. P. (2015). Semiparametric mean field variational Bayes: General principles and numerical issues. *Unpublished*.
- Rice, J. A. and Wu, C. O. (2001). Nonparametric mixed effects models for unequally sampled noisy curves. *Biometrics*, **57**:253–259.
- Richardson, S., Leblond, L., Jaussent, I., and Green, P. J. (2002). Mixture models in measurement error problems, with reference to epidemiological studies. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, **165**:549–566.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo statistical methods*. New York: Springer Science & Business Media.
- Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric regression*, volume 12. New York: Cambridge university press.
- Ruppert, D., Wand, M. P., and Carroll, R. J. (2009). Semiparametric regression during 2003–2007. *Electronic Journal of Statistics*, **3**:1193.
- Ryu, D., Li, E., and Mallick, B. K. (2011). Bayesian nonparametric regression analysis of data with random effects covariates from longitudinal measurements. *Biometrics*,

## REFERENCES

---

- 67**:454–466.
- SAS Institute Inc. (2013). *SAS/STAT 13.1 Users Guide*. North Carolina: SAS Institute Inc.
- Smith, A. D. and Wand, M. P. (2008). Streamlined variance calculations for semiparametric mixed models. *Statistics in Medicine*, **27**:435–448.
- Stan Development Team (2015). *Stan: A C++ Library for probability and sampling. Version 2.6*.
- Stewart, B. (2014). Latent factor regressions for the social sciences. *Unpublished*.
- Tan, L. S. and Nott, D. J. (2013). Variational inference for generalized linear mixed models using partially noncentered parametrizations. *Statistical Science*, **28**:168–188.
- Tan, L. S. and Nott, D. J. (2014). A stochastic variational framework for fitting and diagnosing generalized linear mixed models. *Bayesian Analysis*, **9**:963–1004.
- Teh, Y. W., Newman, D., and Welling, M. (2006). A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation. *Em Advances in neural information processing systems*, pgs. 1353–1360.
- Teschendorff, A. E., Wang, Y., Barbosa-Morais, N. L., Brenton, J. D., and Caldas, C. (2005). A variational Bayesian mixture modelling framework for cluster analysis of gene-expression data. *Bioinformatics*, **21**:3025–3033.
- Titterton, D. (2004). Bayesian methods for neural networks and related models. *Statistical Science*, **23**:128–139.
- Verbyla, A. P., Cullis, B. R., Kenward, M. G., and Welham, S. J. (1999). The analysis of designed experiments and longitudinal data by using smoothing splines. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **48**:269–311.
- Wainwright, M. J. and Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, **1**:1–305.
- Wand, M. P. (2014). Fully simplified multivariate normal updates in non-conjugate variational message passing. *The Journal of Machine Learning Research*, **15**:1351–1369.
- Wand, M. P. (2015). Fast approximate inference for arbitrarily large semiparametric regression models via message passing. *Unpublished*.
- Wand, M. P. and Ormerod, J. T. (2008). On semiparametric regression with O’Sullivan penalized splines. *Australian and New Zealand Journal of Statistics*, **50**:179–198.

## REFERENCES

---

- Wand, M. P. and Ormerod, J. T. (2011). Penalized wavelets: embedding wavelets into semiparametric regression. *Electronic Journal of Statistics*, **5**:1654–1717.
- Wand, M. P., Ormerod, J. T., Padoan, S. A., and Frühwirth, R. (2012). Mean field variational Bayes for elaborate distributions. *Bayesian Analysis*, **7**:847–900.
- Wand, M. P. and Ripley, B. (2006). KernSmooth: Functions for kernel smoothing for wand and jones (1995). *R package version*, **2**:19–22.
- Wang, B. and Titterton, D. (2006). Convergence properties of a general algorithm for calculating variational Bayesian estimates for a normal mixture model. *Bayesian Analysis*, **1**:625–650.
- Wang, C. and Blei, D. M. (2013). Variational inference in nonconjugate models. *The Journal of Machine Learning Research*, **14**:1005–1031.
- Wang, C., Paisley, J. W., and Blei, D. M. (2011). Online variational inference for the hierarchical Dirichlet process. In *International conference on artificial intelligence and statistics*, pgs. 752–760.
- Winkelmann, R. (2004). Health care reform and the number of doctor visits: an econometric analysis. *Journal of Applied Econometrics*, **19**:455–472.
- Winn, J. M. and Bishop, C. M. (2005). Variational message passing. *Journal of Machine Learning Research*, **6**:661–694.
- Wolfinger, R., Tobias, R., and Sall, J. (1994). Computing gaussian likelihoods and their derivatives for general linear mixed models. *SIAM Journal on Scientific Computing*, **14**:1294–1310.
- Wood, S. N. (2014). nlme: Linear and nonlinear mixed effects models. *R package version 3.1*.
- You, C., Ormerod, J. T., and Müller, S. (2014). On variational Bayes estimation and variational information criteria for linear regression models. *Australian and New Zealand Journal of Statistics*, **56**:73–87.
- Zeger, S. L. and Diggle, P. J. (1994). Semiparametric models for longitudinal data with application to CD4 cell numbers in HIV seroconverters. *Biometrics*, pgs. 689–699.
- Zeger, S. L. and Karim, M. R. (1991). Generalized linear models with random effects; a Gibbs sampling approach. *Journal of the American Statistical Association*, **86**:79–86.
- Zeger, S. L. and Liang, K.-Y. (1986). Longitudinal data analysis for discrete and continuous

## REFERENCES

---

- outcomes. *Biometrics*, **42**:121–130.
- Zhang, D., Lin, X., Raz, J., and Sowers, M. (1998). Semiparametric stochastic mixed models for longitudinal data. *Journal of the American Statistical Association*, **93**:710–719.
- Zhao, Y., Staudenmayer, J., Coull, B. A., and Wand, M. P. (2006). General design Bayesian generalized linear mixed models. *Statistical Science*, **21**:35–51.