

Variational Inference for Heteroscedastic and Longitudinal Regression Models

by

Marianne Menictas

*Submitted to the School of Mathematical Sciences, Faculty of Science
in partial fulfilment of the requirements for the degree of*

Doctor of Philosophy

at the

UNIVERSITY OF TECHNOLOGY SYDNEY

2015

© Marianne Menictas, MMXV. All rights reserved.

Permission is herewith granted to UTS 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 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

My deepest thank you goes to my advisor and academic father, Matt. You opened my eyes to a new way of thinking. I could not have asked for a more inspirational, supportive and patient role model to guide me through this journey and will forever be grateful for what you have taught me. I extend my appreciation to my co-supervisor Peter, who spent many hours guiding me with helpful advice throughout my PhD journey. To my honours advisor Debbie, thank you for showing me the importance of rigorous research and for always having your door open to me. Cathy, my PhD buddy and great friend, discussing the intricacies of our work was the highlight of my day so thank you. To Mum and Nick, thank you for your constant unconditional love and support and for always being there for me. Finally, to Dad, my quest would not have taken place if it weren't for you encouraging me to be unreservedly curious and hungry for the pursuit of knowledge. You are the biggest inspiration in my life.

Contents

1	Introduction	1
1.1	Literature review	2
1.1.1	Variational approximations	2
1.1.2	Mean field variational Bayes	3
1.2	Contribution	4
1.3	Notational guide	4
1.3.1	Vector and matrix notation	5
1.3.2	Probability notation	6
1.4	Theorems, definitions and results	7
1.4.1	Vector differential calculus	7
1.4.2	Special functions	8
1.4.3	Probability distributions	8
1.4.4	Useful vector and matrix results	15
1.5	Mean field variational Bayes	16
1.6	Variational message passing	18
1.7	Non-conjugate variational message passing	18
1.8	Graphical models	19
1.8.1	Directed acyclic graphs	20
1.8.2	Moral graphs	23
1.8.3	Graphical models viewpoint of mean field variational Bayes	24
1.9	Accuracy measure of mean field variational Bayes	26
1.10	O’Sullivan penalised splines	26
1.10.1	Mixed model representation	28
1.11	Statistical software	29

1.11.1	Example in BRugs and RStan	29
1.12	The flow of the thesis	31
2	Mean field variational Bayes for marginal longitudinal semiparametric regression	33
2.1	Introduction	33
2.2	Marginal longitudinal semiparametric regression models	34
2.2.1	Additive models extension	36
2.2.2	Incorporation of interactions	38
2.2.3	Varying coefficient models	39
2.3	Mean field variational Bayes methodology	39
2.3.1	Heuristic justification of mean field variational Bayes	44
2.3.2	Mean field variational Bayes approximate density functions for entries of Σ	45
2.4	Simulation study	46
2.4.1	Assessment of speed	48
2.5	Application	50
2.6	Discussion	51
2.A	Derivation of algorithm 1	52
2.B	Derivation for the marginal log-likelihood lower bound	55
2.C	Derivation of the Fisher information matrix for the linear marginal longitudinal model	58
2.D	Proof of Theorem 2.3.1	60
2.E	Proof of Theorem 2.3.2	61
3	Variational inference for heteroscedastic semiparametric regression	62
3.1	Introduction	62
3.2	Model description	63
3.3	Variational inference algorithm	65
3.4	Assessment of performance	69
3.4.1	Assessment of accuracy	69
3.4.2	Assessment of coverage	72
3.4.3	Assessment of speed	74

3.5	Application	74
3.6	Discussion	76
3.A	Derivation of algorithm 2	78
3.B	Derivation for the marginal log-likelihood lower bound	84
3.C	Proof of Proposition 3.A.1	89
4	Heteroscedastic semiparametric regression extensions	91
4.1	Introduction	91
4.2	Bivariate predictor nonparametric regression	92
4.2.1	Model description	93
4.2.2	Simulation study	94
4.2.3	Application	97
4.3	Extension to additive models	98
4.3.1	Application	102
4.4	Real-time heteroscedastic nonparametric regression	104
4.5	Discussion	106
5	Mean field variational Bayes for group-specific curve models	108
5.1	Introduction	108
5.2	Two-level Gaussian response model	109
5.2.1	Mixed model representation	111
5.2.2	Bayesian inference	112
5.2.3	Mean field variational Bayes methodology	113
5.2.4	Streamlining Mean field variational Bayes for the two-level Gaussian response model	117
5.2.5	Simulation study	121
5.2.6	Application to data from a growth study	122
5.3	Three-level Gaussian response model	127
5.3.1	Bayesian three-level Gaussian response model	130
5.3.2	Mean field variational Bayes methodology	131
5.3.3	Streamlining mean field variational Bayes for the three-level Gaus- sian response model	136
5.3.4	Simulation study	141

5.4	Discussion	142
5.A	Derivation of Algorithm 5	143
5.B	Derivation of the marginal log-likelihood lower bound	149
6	Mean field variational Bayes for mixture models in measurement error problems	156
6.1	Introduction	156
6.2	Model structure	157
6.2.1	Nondifferential measurement error	158
6.2.2	Finite normal mixture component	158
6.2.3	Joint model	159
6.3	The full model	160
6.4	Mean field variational Bayes	161
6.5	Simulation study	163
6.6	Discussion	169
6.A	Expectation updates	170
6.B	Derivation of Algorithm 9	171
6.C	Derivation of the marginal log-likelihood lower bound	181
7	Alternative approach based on variational message passing	188
7.1	Introduction	188
7.2	Factor graph representation	189
7.2.1	Additional notation	190
7.3	Natural parameter forms	191
7.4	Primitive integrals and results	195
7.5	Function definitions	197
7.6	The variational message passing algorithm	198
7.7	Examples	200
7.7.1	Marginal longitudinal semiparametric regression model	200
7.7.2	Mixture model with measurement error	205
7.8	Discussion	211
7.A	Derivation of Algorithm 11	212
7.A.1	Step 1: Initialise factor to stochastic node messages	212

CONTENTS

7.A.2	Step 2 (a): Update stochastic node to factor messages	213
7.A.3	Step 2 (b): Update factor to stochastic node messages	214
7.A.4	Step 3 : Optimal q -densities	218
7.B	Derivation of Algorithm 12	219
7.B.1	Step 1: Initialise factor to stochastic node messages	219
7.B.2	Step 2 (a): Update stochastic node to factor messages	222
7.B.3	Step 2 (b): Update factor to stochastic node messages	223
7.B.4	Step 3 : Optimal q -densities	245
8	Conclusion	246

Abstract

The focus of this thesis is on the development and assessment of mean field variational Bayes (MFVB), which is a fast, deterministic tool for inference in a Bayesian hierarchical model setting. We assess the performance of MFVB via the use of comprehensive comparisons against a Markov chain Monte Carlo (MCMC) benchmark. Each of the models considered are special cases of semiparametric regression. In particular, we focus on the development and assessment of the performance of MFVB for heteroscedastic and longitudinal semiparametric regression models. Generally, the new MFVB methodology performs well in its assessment of accuracy against MCMC for the semiparametric and nonparametric regression models considered in this thesis. It is also much faster and is shown to be applicable to real-time analyses. Several real data illustrations are provided. Altogether, MFVB proves to be a credible inference tool and a good alternative to MCMC, especially when analysis is hindered by time constraints.

Chapter 1

Introduction

“Information is the oil of the 21st century & analytics is the combustion engine”

- Peter Sondergaard, Senior Vice President, Gartner Research.

The rapid growth of information is one of the most challenging issues the world faces today. To appreciate the degree by which an information revolution is under way, consider for instance the Sloan Digital Sky Survey telescope in New Mexico, which in 2000, collected more data in its first few weeks of operation than had been amassed in the entire history of astronomy. An even more salient example is Facebook, having just celebrated its 10th birthday, has more than 10 million photos uploaded every hour and a *like* button clicked almost 3 billion times per day. This recent growth of information has become known as the *big data revolution* and acts as the underlying motivation for this thesis. In some instances, the amount of information has grown so much, that the volume being handled no longer fits into standard computer memory. As a result, engineers need to re-develop tools for analysing it all. This is the drive behind new data processing technologies such as Map Reduce and Hadoop, which allows one to manage larger troves of data than ever before.

In a similar vein, the pace of processing information in recent years has given impetus for the development of faster data analysis techniques to replace conventional methods that, however accurate, take much longer to run. Specifically, this thesis focuses on the development and assessment of mean field variational Bayes (MFVB), which is a fast alternative to Markov chain Monte Carlo (MCMC). These methods are considered in a hierarchical Bayesian model setting, with comprehensive comparisons between the two

being a theme of this work. We commence our discussion by considering a review of the literature surrounding MFVB.

1.1 Literature review

The Bayesian statistical literature is dominated by simulation-based approximations such as MCMC. Since its emergence in Bayesian statistics in the 1990s, and its ensuing construction of the Bayesian inference engines BUGS (Bayesian inference Using Gibbs Sampling) (Spiegelhalter *et al.*, 2003) and Stan (Stan Development Team, 2014), MCMC has made previously intractable inference problems solvable. When challenged with a new and possibly complex model that requires fitting and inference, the Bayesian analyst can choose between at least two approaches: (a) to manually program MCMC using Metropolis-Hastings and/or slice sampling as a way of dealing with arduous full conditional distributions; or (b) to implement the model in the BUGS inference engine by the simple press of a button.

Albeit an extremely accurate inference tool, MCMC's main setback from an applied standpoint is the waiting time required for a program to run, sometimes taking days or even weeks to get results. When considering applications that depend on speedy analysis, such as speech recognition and robot vision, these times are unacceptably slow. A rapid deterministic alternative to MCMC and Monte Carlo methods in general is the variational approximation, which involves approximate inference for parameters in complex statistical models. Let us next discuss variational approximations and their emergence into the computer science and statistical literature.

1.1.1 Variational approximations

Variational approximations originate from the mathematical topic of *variational calculus*, which has been part of the mathematical literature for over two centuries. Variational calculus involves optimising a functional over a specified class of functions, somewhat akin to the way in which a function can depend on a numerical variable. Variational approximations arise when constraints are fixed on the specified class of functions, usually to enforce tractability (e.g. Ormerod & Wand, 2010).

The posterior distribution of the model parameters is of primary interest in Bayesian inference. However, even for relatively simple models, the derivation of the posterior

density can sometimes involve intractable integrals. As mentioned previously, our interest lies in the achievement of faster inference for hierarchical Bayesian models. Specifically, we are interested in providing a deterministic alternative to the stochastic MCMC by combining the notions of Bayesian inference together with the variational approximation, in order to approximate these previously intractable integrals. MFVB arises when the variational approximation being used is subject to a specific factorisation, and is therefore, limited in its approximation accuracy.

Variational approximations have made *eye-catching* appearances in the computer science literature, where they have acted as the underlying methodology behind sophisticated areas of research such as Markov random fields, error-control coding and language processing (Jordan *et al.*, 1999; Jordan, 2004). Comprehensive summaries of the variational approximation may be found in Bishop (2006, Chapter 10) and Ormerod & Wand (2010). What is lacking in the computer science literature, however, is research on the accuracy of the variational approximation. This provides the statistical sciences an opportunity to consolidate the existence of an accuracy assessment where we provide a quantitative assessment of the performance of the variational approximation. Research into the quality of the variational approximation has already manifested in statistical literature for specific models. For example, Wang & Titterton (2003) study the consistency properties of variational Bayesian estimators for mixture models involving known densities. It was shown that with probability 1, as the iterations approach infinity, the algorithm converges locally to the maximum likelihood estimator. In addition, Wang *et al.* (2006) introduce a general algorithm for computing variational Bayesian estimates and look into its convergence properties for a normal mixture model. For a comprehensive listing of other relevant literature on the accuracy of the variational approximation, Jordan (2004) and Titterton (2004) discuss further sources.

The specific focus of this thesis therefore, is to provide an in-depth investigation into a specific type of variational approximation, which we refer to as MFVB. The review of the literature surrounding MFVB follows.

1.1.2 Mean field variational Bayes

Mean field variational Bayes (Wainwright & Jordan, 2008) involves Bayesian estimation in graphical models, specifically catered to the variational posterior distribution restricted

to a factorised form. It is one of the most common types of variational approximation used today (Parisi, 1988). The use of MFVB has become more prevalent in the statistical literature in the past decade. For example, Teschendorff *et al.* (2005) use MFVB as a means of accurately categorising tumour types through gene-expression profiling. McGrory & Titterton (2007) applied MFVB to the Bayesian analysis of mixtures of Gaussian distributions. Further, in 2008, a MFVB-based software package named Infer.NET (Minka *et al.*, 2008) was developed. This allowed one to use MFVB on a wide range of statistical problems, making MFVB more accessible than ever before.

In more recent years, MFVB has appeared in the *Journal of the American Statistical Association* (Braun & McAuliffe, 2010; Faes *et al.*, 2011), exposing MFVB to a wider audience and demonstrating its potential application. A new major area of application based on MFVB is real-time semiparametric regression analysis (Luts *et al.*, 2013). Other recent MFVB examples are given in, but not limited to, Hall *et al.* (2011), Huang *et al.* (2013), Wand (2014), McGrory & Titterton (2007), Wang *et al.* (2006) and Neville *et al.* (2014).

1.2 Contribution

Even though we have seen MFVB appear in many different statistical settings, there is still a wide range of research needed to be carried out. This thesis aims to contribute to the *missing parts* of research and application dedicated to MFVB. In particular, our contribution lies in our ability to expose MFVB, yet again, as a faster, more prominent alternative to more widely used inference methods such as the Laplace approximation and MCMC. In addition, we aim to contribute by quantitatively assessing the accuracy of MFVB in various statistical settings.

1.3 Notational guide

There are two main notations used throughout this thesis, namely vector and matrix notation, and probability notation. These are discussed in the following two sections.

1.3.1 Vector and matrix notation

Lower case bold fonts are used to denote vectors and upper case bold fonts denote matrices. The real line is denoted by \mathbb{R} and m -dimensional space as \mathbb{R}^m . The $m \times 1$ column vector with all entries equal to 1 is denoted by $\mathbf{1}_m$. The $m \times m$ identity matrix is denoted by \mathbf{I}_m . For vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$, denoted by $\mathbf{a} = [a_1, \dots, a_m]^\top$ and $\mathbf{b} = [b_1, \dots, b_m]^\top$ we define the following notation.

The element-wise multiplication of two vectors \mathbf{a} and \mathbf{b} is denoted by

$$\mathbf{a} \odot \mathbf{b} = \begin{bmatrix} a_1 b_1 \\ \vdots \\ a_m b_m \end{bmatrix}.$$

The norm of a vector \mathbf{a} is defined to be $\|\mathbf{a}\| = \sqrt{\mathbf{a}^\top \mathbf{a}}$. We use \mathbf{a}_{-i} to denote the vector \mathbf{a} with the i th element removed. A scalar valued function of a vector is evaluated element-wise. For example,

$$\log(\mathbf{a}) = \begin{bmatrix} \log(a_1) \\ \vdots \\ \log(a_m) \end{bmatrix}.$$

We let $\text{diag}(\mathbf{a})$ denote the $m \times m$ diagonal matrix formed by assigning the elements of \mathbf{a} to the main diagonal. For instance,

$$\text{diag}(\mathbf{a}) = \begin{bmatrix} a_1 & 0 & \dots & 0 \\ 0 & a_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_m \end{bmatrix}.$$

For any $m \times m$ matrix \mathbf{A} , $\text{diag}(\mathbf{A})$ refers to the $m \times 1$ vector consisting of the diagonal entries of \mathbf{A} . For square matrices $(\mathbf{A}_1, \dots, \mathbf{A}_d)$, we define $\text{blockdiag}(\mathbf{A}_1, \dots, \mathbf{A}_d)$ to be the block diagonal matrix, with the i th block equal to \mathbf{A}_i . For a symmetric $m \times m$ matrix \mathbf{A} , the determinant and trace are denoted by $|\mathbf{A}|$ and $\text{tr}(\mathbf{A})$ respectively, and are defined in the usual way.

For the matrices $\mathbf{A}(m \times n)$ and $\mathbf{B}(p \times q)$, defined as

$$\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} b_{11} & \vdots & b_{1q} \\ \vdots & \ddots & \vdots \\ b_{p1} & \dots & b_{pq} \end{bmatrix},$$

the following notation is provided.

The transpose of \mathbf{A} is denoted by \mathbf{A}^\top . We let \mathbf{A}^{-1} denote the matrix inverse, given that \mathbf{A} is an invertible matrix. The Kronecker product of two matrices \mathbf{A} and \mathbf{B} is denoted by

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \dots & a_{mn}\mathbf{B} \end{bmatrix}.$$

We define $\text{vec}(\mathbf{A})$ to be the $m^2 \times 1$ vector obtained by stacking the columns of \mathbf{A} underneath each other from left to right. If \mathbf{a} is a $m^2 \times 1$ vector then $\text{vec}^{-1}(\mathbf{a})$ is defined to be the $m \times m$ matrix obtained by listing the entries of \mathbf{a} in a column-wise fashion from left to right. Also, vec^{-1} denotes the usual function inverse since the domain of vec is restricted to square matrices.

1.3.2 Probability notation

If x and y represent scalar random variables, the mean and variance of x are represented as $E(x)$ and $\text{Var}(x)$. The covariance between x and y is represented as $\text{Cov}(x, y)$.

If $\boldsymbol{\theta}$ is a random vector, its density function is denoted by $p(\boldsymbol{\theta})$. If $\boldsymbol{\phi}$ is also a random vector, the conditional density function of $\boldsymbol{\theta}$ given $\boldsymbol{\phi}$ is denoted by $p(\boldsymbol{\theta}|\boldsymbol{\phi})$. In a Bayesian model setting, the full conditional distribution of $\boldsymbol{\theta}$ is given by $p(\boldsymbol{\theta}|\text{rest})$ and represents the conditional distribution of $\boldsymbol{\theta}$ given the rest of the random variables in the model. The mean and covariance of a random vector $\boldsymbol{\theta}$ is given by $E(\boldsymbol{\theta})$ and $\text{Cov}(\boldsymbol{\theta})$.

When using MFVB, the approximate densities are denoted using the letter q . We call these densities q -densities. The mean and variance of a scalar random variable θ under the MFVB criterion is given by:

$$E_q(\theta) = \mu_{q(\theta)} \quad \text{and} \quad \text{Var}_q(\theta) = \sigma_{q(\theta)}^2.$$

Similarly, the mean and covariance of a random vector $\boldsymbol{\theta}$ under the MFVB criterion is given by:

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

Finally, if x_1, \dots, x_n is a sequence of independent random variables following distribution D , we denote that as $x_i \stackrel{\text{ind.}}{\sim} D$, for $1 \leq i \leq n$.

1.4 Theorems, definitions and results

1.4.1 Vector differential calculus

Suppose $f(\mathbf{x})$ is a function with argument $\mathbf{x} \in \mathbb{R}^p$ that returns a scalar. The *derivative vector* of $f(\mathbf{x})$ is the $1 \times p$ vector $D_{\mathbf{x}}f$ with i th entry equal to $\partial f(\mathbf{x})/\partial x_i$ where x_i is the i th entry of \mathbf{x} . $D_{\mathbf{x}}f$ is found using the *First Identification Theorem* (Magnus & Neudecker, 1999):

Theorem 1.4.1. *If \mathbf{a} is a $1 \times p$ vector where $df(\mathbf{x}) = \mathbf{a} d\mathbf{x}$, then the derivative vector is*

$$\mathbf{a} = D_{\mathbf{x}}f.$$

Differentiation again with respect to \mathbf{x} gives the $p \times p$ *Hessian matrix* $H_{\mathbf{x}}f$ with (i, j) th entry equal to $\partial^2 f(\mathbf{x})/\partial x_i \partial x_j$ where x_i and x_j are the i th and j th entry of \mathbf{x} . $H_{\mathbf{x}}f$ is found using the *Second Identification Theorem* (Magnus & Neudecker, 1999):

Theorem 1.4.2. *If \mathbf{A} is a $p \times p$ matrix where $d^2f(\mathbf{x}) = (d\mathbf{x})^\top \mathbf{A} d\mathbf{x}$, then the Hessian matrix is*

$$\mathbf{A} = H_{\mathbf{x}}f.$$

The following results are useful in vector differential calculus.

Result 1.4.1. *If \mathbf{U} and \mathbf{V} are matrix functions, then*

- (a) $\text{vec}(d\mathbf{U}) = d\text{vec}(\mathbf{U})$
- (b) $d|\mathbf{U}| = |\mathbf{U}| \text{tr}(\mathbf{U}^{-1} d\mathbf{U})$
- (c) $d\mathbf{U}^{-1} = -\mathbf{U}^{-1}(d\mathbf{U})\mathbf{U}^{-1}$
- (d) $d(\text{tr}\mathbf{U}) = \text{tr}(d\mathbf{U})$
- (e) $d(\mathbf{U} \otimes \mathbf{V}) = d\mathbf{U} \otimes \mathbf{V} + \mathbf{U} \otimes d\mathbf{V}$
- (f) $d(\log|\mathbf{U}|) = \text{tr}(\mathbf{U}^{-1} d\mathbf{U})$
- (g) $d(\mathbf{U}\mathbf{V}) = (d\mathbf{U})\mathbf{V} + \mathbf{U}(d\mathbf{V})$

1.4.2 Special functions

Definition 1.4.1. *The digamma function is the logarithmic derivative of the gamma function and is defined by*

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

1.4.3 Probability distributions

Here we summarize the main properties of the distributions used throughout this thesis.

1.4.3.1 Uniform distribution

Definition 1.4.2. *A scalar random variable x has a uniform distribution with minimum value a and maximum value b , written $x \sim \text{uniform}(a, b)$ if its density function is*

$$p(x) = \begin{cases} \frac{1}{b-a} & , \quad a \leq x \leq b \\ 0 & , \quad \text{otherwise.} \end{cases}$$

1.4.3.2 Bernoulli distribution

Definition 1.4.3. *A scalar random variable x has a Bernoulli distribution, written $x \sim \text{Bernoulli}(p)$, where $0 \leq p \leq 1$, if its density function is defined to be*

$$p(x) = \begin{cases} p^x(1-p)^{1-x} & , \quad x = 0, 1 \\ 0 & , \quad \text{otherwise.} \end{cases}$$

1.4.3.3 Beta distribution

Definition 1.4.4. A scalar random variable x has a Beta distribution with shape parameters $\alpha, \beta > 0$, written $x \sim \text{Beta}(\alpha, \beta)$, if its density function is defined to be

$$p(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} & , 0 \leq x \leq 1 \\ 0 & , \text{otherwise.} \end{cases}$$

The normalising constant $B(\alpha, \beta)$ is known as the Beta function and is expressed as $\Gamma(\alpha + \beta)/\Gamma(\alpha)\Gamma(\beta)$.

1.4.3.4 Normal distribution

The Normal distribution, also known as the Gaussian distribution, is the most widely used distribution for continuous variables.

Definition 1.4.5. A scalar random variable x has a Normal distribution with mean μ and variance $\sigma^2 > 0$, written $x \sim N(\mu, \sigma^2)$ if its density function is

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}.$$

Definition 1.4.6. The notation $\phi(\cdot)$ is used to denote the standard normal density function, that is $x \sim N(0, 1)$:

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

Thus, for general mean μ and standard deviation σ , that is if $x \sim N(\mu, \sigma^2)$, then

$$p(x) = \frac{1}{\sigma} \phi\left(\frac{x-\mu}{\sigma}\right).$$

Definition 1.4.7. An $n \times 1$ vector of random variables \mathbf{x} has a Multivariate normal distribution with $n \times 1$ mean vector $\boldsymbol{\mu}$ and $n \times n$ positive-definite covariance matrix $\boldsymbol{\Sigma}$, written $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ if its density function is

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

1.4.3.5 Gamma distribution

Definition 1.4.8. A scalar random variable x has a Gamma distribution with shape and rate parameters $A, B > 0$, written $x \sim \text{Gamma}(A, B)$, if its density function is

$$p(x) = \begin{cases} \frac{B^A}{\Gamma(A)} x^{A-1} \exp\{-Bx\} & , x > 0 \\ 0 & , \text{otherwise.} \end{cases}$$

Result 1.4.2. If $x \sim \text{Gamma}(A, B)$, then for any $c > 0$,

$$cx \sim \text{Gamma}(A, B/c).$$

1.4.3.6 Inverse-Gamma distribution

Perhaps the main use of the Inverse-Gamma distribution is in Bayesian statistics, where it is used as an analytically tractable prior distribution for the variance parameters.

Definition 1.4.9. A scalar random variable x has an Inverse Gamma distribution with shape and rate parameters $A, B > 0$, written $x \sim \text{Inverse-Gamma}(A, B)$ if its density function is

$$p(x) = \begin{cases} \frac{B^A}{\Gamma(A)} x^{-A-1} \exp\left\{-\frac{B}{x}\right\} & , x > 0 \\ 0 & , \text{otherwise.} \end{cases}$$

Result 1.4.3. *If $x \sim \text{Inverse-Gamma}(A, B)$, then*

$$E(1/x) = A/B, \quad E\{\log(x)\} = \log(B) - \psi(A),$$

where ψ is the digamma function as defined in Definition 1.4.1

Result 1.4.4. *If $x \sim \text{Gamma}(A, B)$, then*

$$1/x \sim \text{Inverse-Gamma}(A, B).$$

Result 1.4.5. *If $x \sim \text{Gamma}(\frac{\nu}{2}, \frac{1}{2})$, then x is identical to the chi-squared distribution with ν degrees of freedom χ_ν^2 .*

1.4.3.7 Half-Cauchy distribution

Often, alternative scale parameter priors are desired. Gelman (2006) argues that Half-Cauchy densities are better for achieving non-informativeness of scale parameters. The Half-Cauchy distribution is a special case of the conditionally-conjugate folded-noncentral t family of prior distributions for the standard deviation parameter and tends to give cleaner inferences.

Definition 1.4.10. *A scalar random variable x has a Half-Cauchy distribution with scale parameter $A > 0$, written $x \sim \text{Half-Cauchy}(A)$, if its density function is*

$$p(x) = \begin{cases} \frac{2}{\pi A \{1+(x/A)^2\}} & , \quad x > 0 \\ 0 & , \quad \text{otherwise.} \end{cases}$$

The following result is Result 5 of Wand, Ormerod, Padoan and Frühwirth (2011):

Result 1.4.6. Suppose that x and a are random variables such that

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

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

1.4.3.8 Wishart distribution

Definition 1.4.11. An $n \times n$ matrix \mathbf{X} has a Wishart distribution with degrees of freedom $a > 0$ and positive definite rate matrix \mathbf{B} , written $\mathbf{X} \sim \text{Wishart}(a, \mathbf{B})$, if its density function is

$$p(\mathbf{X}) = \begin{cases} C_{n,a}^{-1} |\mathbf{B}|^{a/2} |\mathbf{X}|^{(a-n-1)/2} \exp\{-\frac{1}{2}\text{tr}(\mathbf{B}\mathbf{X})\} & , \quad \mathbf{X} \text{ symmetric} \\ & \text{and positive definite} \\ 0 & , \quad \text{otherwise,} \end{cases}$$

where $C_{n,a} \equiv 2^{an/2} \pi^{n(n-1)/4} \prod_{i=1}^n \Gamma((a+1-i/2))$.

Result 1.4.7. If $\mathbf{X} \sim \text{Wishart}(a, \mathbf{B})$, then

$$E(\mathbf{X}) = a\mathbf{B}^{-1}.$$

1.4.3.9 Inverse-Wishart distribution

The Inverse-Wishart distribution is commonly used as the conjugate prior for the covariance matrix of a multivariate normal distribution.

Definition 1.4.12. An $n \times n$ matrix \mathbf{X} has an Inverse-Wishart distribution with degrees of freedom $a > 0$ and positive definite scale matrix \mathbf{B} , written $\mathbf{X} \sim \text{Inverse-Wishart}(a, \mathbf{B})$, if its density function is

$$p(\mathbf{X}) = \begin{cases} C_{n,a}^{-1} |\mathbf{B}|^{a/2} |\mathbf{X}|^{-(a+n+1)/2} \exp\{-\frac{1}{2}\text{tr}(\mathbf{B}\mathbf{X}^{-1})\} & , \mathbf{X} \text{ symmetric} \\ & \text{and positive definite} \\ 0 & , \text{otherwise,} \end{cases}$$

where $C_{n,a} \equiv 2^{an/2} \pi^{n(n-1)/4} \prod_{i=1}^n \Gamma((a+1-i/2))$.

Result 1.4.8. If $\mathbf{X} \sim \text{Wishart}(a, \mathbf{B})$, then

$$\mathbf{X}^{-1} \sim \text{Inverse-Wishart}(a, \mathbf{B}).$$

1.4.3.10 Multinomial distribution

Definition 1.4.13. A vector of random variables $\mathbf{x} = (x_1, \dots, x_k)$ where x_i takes on values $0, \dots, n$, has a Multinomial distribution with number of independent trials $n > 0$ and probabilities of success $\mathbf{p} = (p_1, \dots, p_k)$, written $x_1, \dots, x_k \sim \text{Multinomial}(n, p_1, \dots, p_k)$, if its density function is

$$p(x_1, \dots, x_k) = \begin{cases} \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}, & \text{when } \sum_{i=1}^k x_i = n \\ 0, & \text{otherwise} \end{cases}.$$

Result 1.4.9. Suppose that $x_1, \dots, x_k \sim \text{Multinomial}(n, p_1, \dots, p_k)$. Then,

$$E(x_i) = np_i.$$

1.4.3.11 Dirichlet distribution

Definition 1.4.14. A vector of random variables $\mathbf{x} = (x_1, \dots, x_k)$ has a Dirichlet distribution with parameters $\alpha_1, \dots, \alpha_k > 0$ written $p(x_1, \dots, x_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ if its density function is

$$p(x_1, \dots, x_k) = \frac{1}{B(\alpha_1, \dots, \alpha_k)} \prod_{i=1}^k x_i^{\alpha_i-1}$$

for all $x_1, \dots, x_k > 0$ satisfying $\sum_{i=1}^k x_i = 1$. The normalizing constant is known as the multinomial Beta function, which can be expressed in terms of the Gamma function as:

$$B(\alpha_1, \dots, \alpha_k) = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)}.$$

Result 1.4.10. Suppose that $x_1, \dots, x_k \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$. Then,

$$E(x_i) = \alpha_i / \sum_{k=1}^K \alpha_k, \quad E\{\log(x_i)\} = \psi(\alpha_i) - \psi\left(\sum_{k=1}^K \alpha_k\right).$$

1.4.3.12 Normal mixture distribution

Definition 1.4.15. A scalar random variable x has a normal mixture distribution with K weighted sum components, written $p(x; \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \sim \text{Normal-Mixture}(\boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\sigma}^2)$, if its density function is

$$p(x; \boldsymbol{\omega}, \boldsymbol{\mu}, \boldsymbol{\sigma}^2) = \sum_{k=1}^K \omega_k (2\pi\sigma_k^2)^{-1/2} \exp\left\{-\frac{1}{2}(x - \mu_k)^2/\sigma_k^2\right\},$$

where $\boldsymbol{\omega} = (\omega_1, \dots, \omega_K)$ are the mixture weights, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$ are the mixture means and $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_K^2)$ are the mixture variances associated with each k th component density.

1.4.4 Useful vector and matrix results

Result 1.4.11. *If \mathbf{a} and \mathbf{b} are vectors of equal length then $\text{diag}(\mathbf{a})\mathbf{b} = \mathbf{a} \odot \mathbf{b}$.*

Result 1.4.12. *If \mathbf{A} and \mathbf{B} are matrices such that the product \mathbf{AB} can be formed, then*

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}).$$

Result 1.4.13. *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 1.4.14. *If \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} are matrices such that the products \mathbf{AC} and \mathbf{BD} can be formed, then*

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}).$$

Result 1.4.15. *If \mathbf{A} is an $m \times m$ matrix and \mathbf{B} is an $n \times n$ matrix, then*

$$\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A})\text{tr}(\mathbf{B}).$$

Result 1.4.16. *If \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} are matrices such that \mathbf{ABCD} is a square matrix, then*

$$\text{tr}(\mathbf{ABCD}) = \text{vec}(\mathbf{D})^\top (\mathbf{A} \otimes \mathbf{C}^\top) \text{vec}(\mathbf{B}^\top).$$

Result 1.4.17. *Let \mathbf{x} be a random vector. Then*

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

Result 1.4.18. *Let \mathbf{x} be a random vector and \mathbf{A} be a random matrix. Then*

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

Result 1.4.19. *Let \mathbf{a} be constant vector, let \mathbf{b} be a random vector, and let \mathbf{C} be a constant matrix with the same number of rows as \mathbf{b} . Then*

$$E(\|\mathbf{a} - \mathbf{C}\mathbf{b}\|^2) = \|\mathbf{a} - \mathbf{C}E(\mathbf{b})\|^2 + \text{tr} \{\mathbf{C} Cov(\mathbf{b}) \mathbf{C}^\top\}.$$

Result 1.4.20. *Let \mathbf{x} and \mathbf{y} be $n \times 1$ random vectors such that \mathbf{x} is conditioned on \mathbf{y} . Then,*

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

1.5 Mean field variational Bayes

Let us consider a general Bayesian model with parameter vector $\boldsymbol{\theta}$ and observed data vector \mathbf{y} . As discussed earlier, even for simple models, it is often the case that the posterior density function $p(\boldsymbol{\theta}|\mathbf{y})$ is intractable. To overcome this problem, we employ MFVB where the objective is to find a tractable distribution $q(\boldsymbol{\theta})$ that closely approximates the true posterior distribution $p(\boldsymbol{\theta}|\mathbf{y})$. For example, we wish to replace the joint posterior density function $p(\boldsymbol{\theta}|\mathbf{y})$ with an approximating density function $q(\boldsymbol{\theta})$ that assumes the product density form

$$q(\boldsymbol{\theta}) = \prod_{i=1}^k q_i(\boldsymbol{\theta}_i) \tag{1.1}$$

for some partition $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k\}$ of $\boldsymbol{\theta}$. Tractability is a major factor on the choice of partition of the approximate density function $q(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$. The so-called q -densities are then chosen to minimise the Kullback-Leibler distance between the two density functions:

$$\int \log \left\{ \frac{p(\boldsymbol{\theta}|\mathbf{y})}{q(\boldsymbol{\theta})} \right\} q(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

subject to (1.1). Note that the integrals are replaced with sums in the case of discrete variables, however in this thesis we assume $\boldsymbol{\theta}$ to be continuous over the parameter space Θ . An equivalent optimization problem, attained through standard manipulations, is to maximise:

$$\underline{p}(\mathbf{y}; q) \equiv \exp \int q_1(\boldsymbol{\theta}_1) \dots q_k(\boldsymbol{\theta}_k) \log \left\{ \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q_1(\boldsymbol{\theta}_1) \dots q_k(\boldsymbol{\theta}_k)} \right\} d\boldsymbol{\theta}_1 \dots d\boldsymbol{\theta}_k \quad (1.2)$$

where $\underline{p}(\mathbf{y}; q)$ represents the lower bound on the marginal likelihood $p(\mathbf{y})$ for all q -densities. Taking the logarithm of (1.2) gives

$$\log \underline{p}(\mathbf{y}; q) \propto \int q_1(\boldsymbol{\theta}_1) \log \tilde{p}(\mathbf{y}, \boldsymbol{\theta}_1) d\boldsymbol{\theta}_1 - \int q_1(\boldsymbol{\theta}_1) \log q_1(\boldsymbol{\theta}_1) d\boldsymbol{\theta}_1 \quad (1.3)$$

where

$$\begin{aligned} \log \tilde{p}(\mathbf{y}, \boldsymbol{\theta}_1) &= \int \log p(\mathbf{y}, \boldsymbol{\theta}) q_2(\boldsymbol{\theta}_2) \times \dots \times q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_2 \dots d\boldsymbol{\theta}_k \\ &= E_{-\boldsymbol{\theta}_1} \{ \log p(\mathbf{y}, \boldsymbol{\theta}) \}, \end{aligned}$$

where $E_{-\boldsymbol{\theta}_1}$ is the expectation with respect to all factors except $q_1(\boldsymbol{\theta}_1)$. By recognising that (1.3) is the negative Kullback Leibler divergence between $q_1(\boldsymbol{\theta}_1)$ and $\tilde{p}(\mathbf{y}, \boldsymbol{\theta}_1)$, we note that the minimum occurs when $q_1(\boldsymbol{\theta}_1) = \tilde{p}(\mathbf{y}, \boldsymbol{\theta}_1)$. Thus we obtain a general expression for the optimal solution $q_1^*(\boldsymbol{\theta}_1)$ given by

$$q_1^*(\boldsymbol{\theta}_1) \propto \exp \{ E_{-\boldsymbol{\theta}_1} \log p(\mathbf{y}, \boldsymbol{\theta}) \}.$$

This can be generalized to an expression for the optimal q -density for each of the i factors:

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

It is easy to show that an alternative expression for $q_i^*(\boldsymbol{\theta}_i)$ is

$$q_i^*(\boldsymbol{\theta}_i) \propto \exp \{ E_{-\boldsymbol{\theta}_i} \log p(\boldsymbol{\theta}_i | \text{rest}) \}, \quad 1 \leq i \leq k, \quad (1.4)$$

where ‘rest’ is the set comprising the random variables in the model, except for $\boldsymbol{\theta}_i$. An iterative coordinate ascent scheme induced by (1.4) is then used to find the optimal parameters in these q -densities. If conjugate priors have been used, then the optimal q -densities belong to known density families (e.g., Winn & Bishop, 2005).

Assuming that each iteration results in an increase in $\log \underline{p}(\mathbf{y}; q)$, and the search restricted to a compact set, then convergence is guaranteed to a local maximiser of $\log \underline{p}(\mathbf{y}; q)$

(Luenberger & Ye, 2008, p. 253). Convergence can be monitored using successive values of $\log \underline{p}(\mathbf{y}; q)$.

1.6 Variational message passing

Variational message passing (VMP) is an alternative to MFVB which results in exactly the same answers, but instead works with “messages” that are passed between neighbouring nodes on the factor graph of the model. Factor graphs will be discussed in detail in Chapter 7. VMP is primarily concerned with conjugate exponential family distributions, represented by their natural parameter versions of the distribution. An in-depth explanation as well as a discussion on the advantages and disadvantages of VMP is given in Chapter 7.

1.7 Non-conjugate variational message passing

Non-conjugate variational message passing (NCVMP) is an extension of VMP which obviates the restriction of only being able to work with conjugate exponential family distributions, whilst maintaining the purely algebraic form of the optimal parameters (Knowles & Minka, 2011). The modularity of NCVMP allows modification only to those model parameters that involve previously intractable solutions, whilst keeping the VMP or MFVB solutions of the other parameters in the model. As a result, NCVMP broadens the class of tractable models for approximate inference using VMP and/or MFVB.

Suppose we approximate the joint posterior density function $p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathbf{y})$ by the product density form

$$q_1(\boldsymbol{\theta}_1) \dots q_k(\boldsymbol{\theta}_k) q_\phi(\boldsymbol{\phi}), \quad (1.5)$$

where $\boldsymbol{\theta} = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k\}$. From Section 1.5 we see that the solutions satisfy

$$\begin{aligned} q_i^*(\boldsymbol{\theta}_i) &\propto \exp \{E_{q(-\boldsymbol{\theta}_i)} \log p(\boldsymbol{\theta}_i | \mathbf{y}, \boldsymbol{\theta} \setminus \boldsymbol{\theta}_i, \boldsymbol{\phi})\}, \quad 1 \leq i \leq k, \\ q_\phi^*(\boldsymbol{\phi}) &\propto \exp \{E_{q(-\boldsymbol{\phi})} \log p(\boldsymbol{\phi} | \mathbf{y}, \boldsymbol{\theta})\}, \end{aligned}$$

where $\boldsymbol{\theta} \setminus \boldsymbol{\theta}_i$ represents the model parameters $\boldsymbol{\theta}$ without $\boldsymbol{\theta}_i$. When $E_{q(-\boldsymbol{\phi})} \log p(\boldsymbol{\phi} | \mathbf{y}, \boldsymbol{\theta})$ does not offer a tractable solution, NCVMP can be used by replacing (1.5) with

$$q_1(\boldsymbol{\theta}_1) \dots q_k(\boldsymbol{\theta}_k) q_\phi(\boldsymbol{\phi}; \boldsymbol{\eta}),$$

where $q_\phi(\boldsymbol{\phi}; \boldsymbol{\eta})$ is forced to have an exponential family density function with natural parameter vector $\boldsymbol{\eta}$ and natural statistic $\mathbf{T}(\boldsymbol{\phi})$. Then, using Theorem 1 of Knowles & Minka (2011) to find an update for $\boldsymbol{\eta}$ and regular MFVB for the other parameters, we get

$$\begin{aligned} q_i^*(\boldsymbol{\theta}_i) &\propto \exp \{E_{q(-\boldsymbol{\theta}_i)} \log p(\boldsymbol{\theta}_i | \mathbf{y}, \boldsymbol{\theta} \setminus \boldsymbol{\theta}_i, \boldsymbol{\phi})\}, \quad 1 \leq i \leq k, \\ \boldsymbol{\eta} &\leftarrow [\text{var} \{\mathbf{T}(\boldsymbol{\phi})\}]^{-1} [\mathbf{D}_\eta E_{q(\boldsymbol{\theta}, \boldsymbol{\phi})} \{\log p(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\phi})\}]^\top. \end{aligned} \quad (1.6)$$

This thesis only considers the case where $q(\boldsymbol{\phi}; \boldsymbol{\eta})$ corresponds to a d -dimensional Multivariate Normal density function. For this special case, Wand (2014) explains that the number of entries in $\text{var} \{\mathbf{T}(\boldsymbol{\phi})\}$ is quartic in d and thus the $\boldsymbol{\eta}$ update in (1.6) becomes numerically challenging. To get around this, Wand (2014) proves that

$$\boldsymbol{\eta} \leftarrow [\text{var} \{\mathbf{T}(\boldsymbol{\phi})\}]^{-1} [\mathbf{D}_\eta E_{q(\boldsymbol{\theta}, \boldsymbol{\phi})} \{\log p(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\phi})\}]^\top$$

is mathematically equivalent to the following updates

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

where

$$S = E_{q(\boldsymbol{\theta}, \boldsymbol{\phi})} \{\log p(\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\phi})\}.$$

vec and vec^{-1} are defined in Section 1.3.1 and \mathbf{D} represents the derivative vector which is defined in Section 1.4.1. However, to get around using the vec operator, a result given in the appendix of Opper & Archambeau (2009) shows that an equivalent representation of the second of these updates is

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\phi})} \leftarrow \left(-\mathbf{H}_{\boldsymbol{\mu}_{q(\boldsymbol{\phi})}} S \right)^{-1}$$

where \mathbf{H} represents the Hessian matrix and is defined in Section 1.4.1. The proof of this result is given in Appendix 3.C of this thesis. We work with this alternative form when using NCVMP.

1.8 Graphical models

One of the engaging aspects of graphical models is that a particular graph has the ability to make probabilistic statements for an extensive class of distributions. In particular, we focus

on the use of graph representations of hierarchical Bayesian models, as they are an effective tool in understanding the conditional dependence structure of these models. Graphical models aid in the calculation of the imposed product factorisation on the posterior $p(\boldsymbol{\theta}|\mathbf{y})$ which is involved in the initial stages of MFVB. As a result, an additional partition of the dependence structure in the model, known as an induced factorization (Bishop, 2006), is then possible. Both of these factorizations are subject to the type of imposed factorization and the underlying structure of the model. Graphical models are a very useful tool in aiding our understanding of these factorizations. These are defined and outlined in more detail in the following section.

1.8.1 Directed acyclic graphs

A graph comprises a set of nodes together with a set of edges, where edges are placed between pairs of nodes. In this thesis, nodes are depicted as circles and edges are depicted as line segments. Also, we use the notation $a \perp\!\!\!\perp b|c$ to denote that a is conditionally independent to b given c , in a graphical model setting.

Definition 1.8.1. *A graph is an **undirected graph** if it comprises a set of nodes connected by undirected edges.*

Definition 1.8.2. *A graph is a **directed graph** if it comprises a set of nodes connected by directed edges.*

Definition 1.8.3. *In a directed graph, a path between a set of nodes is a **cycle** if all directed edges in the path meet the nodes head-to-tail throughout the path.*

Figure 1.1 illustrates examples of an undirected graph, directed graph and a cycle.

Definition 1.8.4. *A directed graph which has no cycles is **acyclic**. This is referred to as a **directed acyclic graph (DAG)**.*

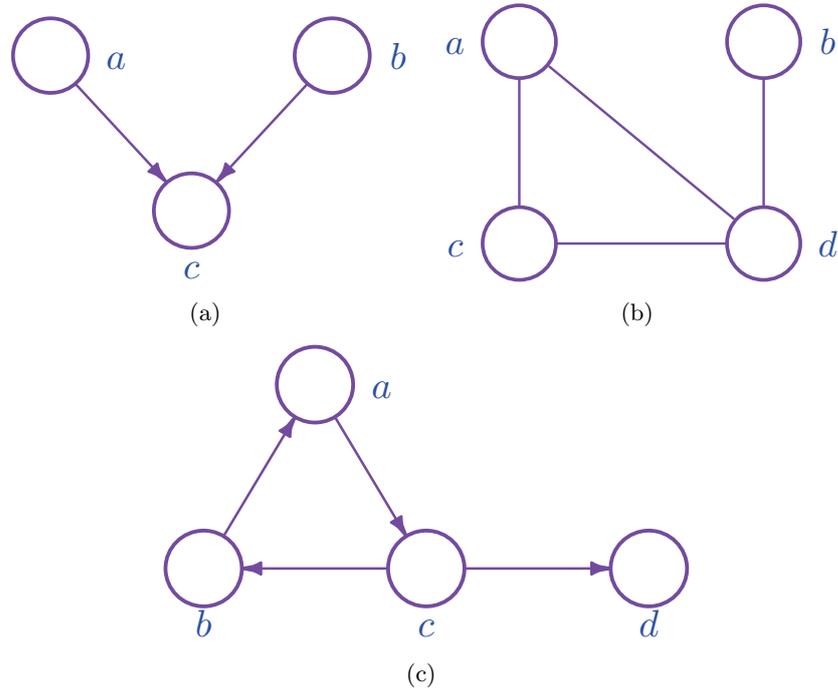


Figure 1.1: (a) Example of a directed graph. (b) Example of an undirected graph. (c) A directed graph which has cycle (a, c, b) .

Throughout this thesis, we will use *DAG* as shorthand for *directed acyclic graph*. Figure 1.2 displays a simple example of a DAG.

As mentioned previously, graphical representations of hierarchical Bayesian models are an effective tool in understanding the conditional dependence structure of such models. This is achieved when we consider the nodes on a DAG to represent the random variables in the model, and the directed edges to represent the conditional dependence structure of the model. It is often easier to think of a DAG as a *family tree*, where a directed edge represents a *child-parent* relationship within the family.

Definition 1.8.5. In a DAG, if two nodes are connected by a directed edge then the node connected by the arrow-head is the **child** and the node connected by the tail is the **parent**. If two nodes have a child in common, they are known as **co-parents**.

Definition 1.8.6. Within a DAG, the **Markov blanket** of a node is the set containing the node's children, parents and co-parents.

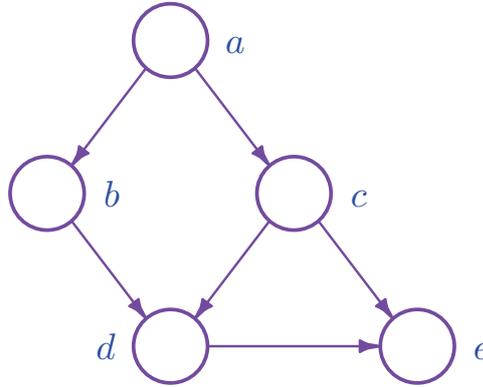
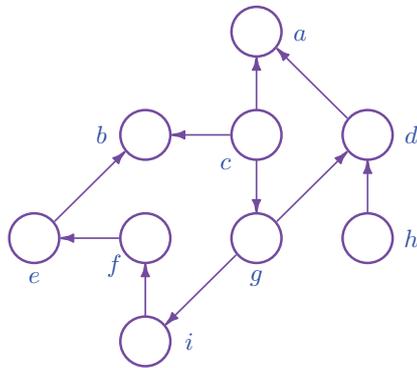


Figure 1.2: A simple DAG containing nodes a, b, c, d and e .



node	children	parents	co-parents	Markov blanket
a	$\{\}$	$\{c, d\}$	$\{\}$	$\{c, d\}$
b	$\{\}$	$\{c, e\}$	$\{\}$	$\{c, e\}$
c	$\{a, b, g\}$	$\{\}$	$\{d, e\}$	$\{a, b, d, e, g\}$
d	$\{a\}$	$\{g, h\}$	$\{c\}$	$\{a, c, g, h\}$
e	$\{b\}$	$\{f\}$	$\{c\}$	$\{b, c, f\}$
f	$\{e\}$	$\{i\}$	$\{\}$	$\{e, i\}$
g	$\{d, i\}$	$\{c\}$	$\{h\}$	$\{c, d, h, i\}$
h	$\{d\}$	$\{\}$	$\{g\}$	$\{d, g\}$
i	$\{f\}$	$\{g\}$	$\{\}$	$\{f, g\}$

Figure 1.3: Left panel: Example of a DAG. Right panel: Table of each node's corresponding children, parents, co-parents and Markov blanket from the given DAG.

Figure 1.3 illustrates the *family tree* terminology for a given DAG, including its Markov blanket. We now turn our attention to probabilistic DAGs where the nodes represent random variables $\{x_1, \dots, x_k\}$ with joint density function $p(x_1, \dots, x_k)$, given by Definition 1.8.7.

Definition 1.8.7. The joint density function of the random variables x_1, \dots, x_k represented by nodes in a DAG is given by

$$p(x_1, \dots, x_k) = \prod_{i=1}^k p(x_i \mid \text{parents of } x_i).$$

DAGs aid in the calculation of the full conditional distributions as is clear from the following definition.

Definition 1.8.8. *The full conditional distribution of a node x_i is the distribution of x_i given all of the other random variables in the DAG, and is denoted by $p(x_i|rest)$.*

The following result shows that probabilistically we can separate each node and its corresponding Markov blanket with the remainder of the DAG.

Result 1.8.1. *Let $p(x_1, \dots, x_k)$ be a density function defined on a DAG. Then*

$$p(x_i|rest) = p(x_i|\text{Markov blanket of } x_i), \quad 1 \leq i \leq k.$$

Result 1.8.1 corresponds to Section 3.2 of Jordan (2004). Making use of Result 1.8.1 in equation (1.4) we now have the optimal q -densities satisfying

$$q_i^*(x_i) \propto \exp\{E_{-x_i} \log p(x_i|\text{Markov blanket of } x_i)\}, \quad 1 \leq i \leq k. \quad (1.7)$$

This is known as the *locality property* of DAGs and can be extremely important when performing MFVB on large models.

1.8.2 Moral graphs

Moral graphs are undirected graphs that are constructed from DAGs by forcing the *im-moral* parents of a common node to get *married*.

Definition 1.8.9. *The **moral graph** of a DAG is constructed by adding an undirected edge between all pairs of parents of a node that are not already connected and turning all directed edges into undirected edges. This procedure is referred to as **moralisation**.*

The following two definitions deal with ancestral sets, which are useful for results concerning probabilistic graphs.

Definition 1.8.10. An *ancestral set* is a set of nodes in a DAG such that for each node in the set, all of its parents are also in the set.

Definition 1.8.11. Let S be a subset of nodes in a DAG. Then the ancestral set containing S that has the fewest number of nodes is the **smallest ancestral set containing S** .

The above definitions aid in the construction of the Theorem 1.8.1.

Theorem 1.8.1. If A , B and C are disjoint node subsets in a DAG then $A \perp\!\!\!\perp B|C$ if C separates A from B in the moral graph of the smallest ancestral set containing $A \cup B \cup C$.

Theorem 1.8.1 corresponds to Corollary 3.23 of Lauritzen (1996).

1.8.3 Graphical models viewpoint of mean field variational Bayes

When considering the conditional independence structure of a DAG, the main benefit of incorporating moral graphs and ancestral sets for MFVB is their extremely simple nature compared to that of d-separation (Bishop, 2006). For example, consider the Bayesian hierarchical Gaussian linear mixed model:

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2\mathbf{I}), & \mathbf{u}|\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2 &\sim \mathbf{N}(\mathbf{0}, \text{blockdiag}(\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2)), \\ \boldsymbol{\beta} &\sim \mathbf{N}(\mathbf{0}, \sigma_\beta^2\mathbf{I}), & \sigma_{u_\ell}^2 &\sim \text{Inverse-Gamma} \sim (A_{u_\ell}, B_{u_\ell}), \quad 1 \leq \ell \leq 3, \\ & & \sigma_\varepsilon^2 &\sim \text{Inverse-Gamma} \sim (A_\varepsilon, B_\varepsilon), \end{aligned} \tag{1.8}$$

for $\sigma_\beta^2, A_{u_\ell}, B_{u_\ell}, A_\varepsilon, B_\varepsilon > 0$ and where \mathbf{y} is a vector of responses and \mathbf{X} and \mathbf{Z} are design matrices. Also, $\boldsymbol{\beta}$ and \mathbf{u} are vectors of fixed effects and random effects, respectively, and $\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2$ and σ_ε^2 are variance parameters. Figure 1.4 shows the DAG and corresponding moral graph for (1.8) where the shaded node corresponds to the observed data vector. Random effects and auxiliary variables are referred to as hidden nodes. The conditional

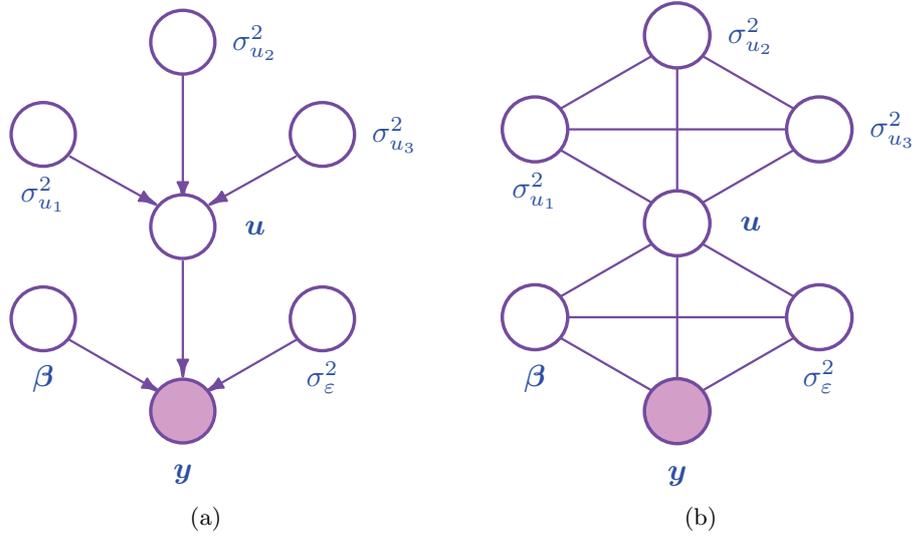


Figure 1.4: (a) DAG representing model (1.8). (b) Moral graph of DAG in (a).

dependence structure of (1.8) is illustrated through the directed edges in the DAG in Figure 1.4 (a).

MFVB is focused around imposing a product restriction on the joint posterior density function. For instance, we may impose the restriction

$$q(\boldsymbol{\beta}, \mathbf{u}, \sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2, \sigma_\varepsilon^2) = q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2, \sigma_\varepsilon^2). \quad (1.9)$$

The right-hand side of (1.9) corresponds to removing the edges connecting $\boldsymbol{\beta}$ with σ_ε^2 , \mathbf{u} with σ_ε^2 , \mathbf{u} with $\sigma_{u_1}^2$, \mathbf{u} with $\sigma_{u_2}^2$, and \mathbf{u} with $\sigma_{u_3}^2$ in the moral graph in Figure 1.4 (b). A tractable solution indeed arises for this product restriction, however to illustrate the simplicity of using moralisation for induced factorisations, we assume that a further breakdown is necessary.

Let $A = \{\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2\}$, $B = \sigma_\varepsilon^2$ and $C = \{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\}$. The smallest ancestral set containing $A \cup B \cup C$ is the entire DAG. It is also evident from the moral graph that all paths between $\{\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2\}$ and σ_ε^2 must pass through at least one of $\{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\}$. Therefore, adopting Theorem 1.8.1 we see that

$$\{\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2\} \perp\!\!\!\perp \sigma_\varepsilon^2 \mid \{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\},$$

and hence (1.9) is updated to

$$q(\boldsymbol{\beta}, \mathbf{u}, \sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2, \sigma_\varepsilon^2) = q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma_{u_1}^2, \sigma_{u_2}^2, \sigma_{u_3}^2) q(\sigma_\varepsilon^2).$$

Often, the use of induced factorisations is very beneficial, especially when considering increasingly complex models. The use of moral graphs enables detection of these factorisations which results in efficient derivations of MFVB solutions.

1.9 Accuracy measure of mean field variational Bayes

The accuracy of MFVB can range from excellent to relatively poor, depending on the structure of the DAG and type of product restriction imposed on the model i.e. number of edges removed. This section gives detail on the accuracy measure used to assess the performance of MFVB against its MCMC benchmark.

Let θ represent a generic univariate parameter. We are interested in measuring the accuracy of an MFVB approximate density $q^*(\theta)$ with respect to the exact posterior density $p(\theta|\mathbf{y})$. There are various ways in which this can be done. Faes *et al.* (2011) recommend working with the L_1 loss, or integrated absolute error (IAE) of $q^*(\theta)$. This quantity is given by

$$\text{IAE} \{q^*(\theta)\} = \int_{-\infty}^{\infty} |q^*(\theta) - p(\theta|\mathbf{y})| d\theta.$$

The fact that $\text{IAE} \{q^*(\theta)\}$ lies between 0 and 2 gives motivation to the accuracy measure

$$\text{accuracy} \{q^*(\theta)\} = 1 - \frac{1}{2} \text{IAE} \{q^*(\theta)\},$$

where the quantity $\frac{1}{2} \text{IAE} \{q^*(\theta)\}$ is the total variation metric corresponding to the probability measure induced by $p(\theta|\mathbf{y})$ and $q^*(\theta)$. Since $0 \leq \text{accuracy} \{q^*(\theta)\} \leq 1$ and thus can be expressed as a percentage. The posterior $p(\theta|\mathbf{y})$ can be approximated quite well by performing MCMC with sufficiently large sample sizes.

1.10 O’Sullivan penalised splines

O’Sullivan splines are an immediate generalisation of smoothing splines (e.g. Green & Silverman, 1994). Prior to the mid-1990s, smoothing splines were of particular interest in the spline-based non-parametric regression literature. This involved the number of basis functions approximately equal to the sample size (e.g. Wahba, 1990; Green & Silverman, 1994). O’Sullivan penalised splines possess the attractive feature of using considerably fewer basis functions, and their natural boundary conditions make them a very reliable choice of basis

in semiparametric and nonparametric regression. Wand & Ormerod (2008) give a detailed description of how O'Sullivan penalised splines can be added to the mixed model-based regression setting including examples with R code. We summarise the approach here.

Consider the simple 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 an estimate of f over the interval $[a, b]$ comprising the x_i s. For $K \leq n$, let B_1, \dots, B_{K+4} be the cubic-B-spline basis functions defined by the knot sequence $a = \kappa_1 = \kappa_2 = \kappa_3 = \kappa_4 < \kappa_5 < \dots < \kappa_{K+4} < \kappa_{K+5} = \kappa_{K+6} = \kappa_{K+7} = \kappa_{K+8} = b$. We require a function that minimises the penalised residual sum of squares

$$\text{RSS}(f, \lambda) = \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda \int_a^b \{f''(x)\}^2 dx, \quad (1.10)$$

where $\lambda > 0$ is a smoothing parameter. The solution is the O'Sullivan penalised spline $f(\mathbf{x}) = \mathbf{B}\boldsymbol{\nu}$ and thus (1.10) reduces to

$$\text{RSS}(\boldsymbol{\nu}, \lambda) = (\mathbf{y} - \mathbf{B}\boldsymbol{\nu})^\top (\mathbf{y} - \mathbf{B}\boldsymbol{\nu}) + \lambda \boldsymbol{\nu}^\top \boldsymbol{\Omega} \boldsymbol{\nu} \quad (1.11)$$

where $B_{ik} = B_k(x_i)$, and $\Omega_{kk'} = \int_a^b B_k''(x) B_{k'}''(x) dx$. It is easy to see that the solution to (1.11) is

$$\hat{\boldsymbol{\nu}} = (\mathbf{B}^\top \mathbf{B} + \lambda \boldsymbol{\Omega})^{-1} \mathbf{B}^\top \mathbf{y} \quad (1.12)$$

and the fitted O'Sullivan penalised spline is thus given by $\hat{f}(\mathbf{x}) = \mathbf{B}\hat{\boldsymbol{\nu}}$. Computation of the design matrix \mathbf{B} is fairly straightforward, however computation of $\boldsymbol{\Omega}$ can be somewhat challenging. Using the theorem in Section 6 of Wand & Ormerod (2008) we get an expression for $\boldsymbol{\Omega}$:

$$\boldsymbol{\Omega} = (\tilde{\mathbf{B}}'')^\top \text{diag}(\boldsymbol{\omega}) \tilde{\mathbf{B}}'' \quad (1.13)$$

where $\tilde{\mathbf{B}}''$ is the $3(K+7) \times (K+4)$ matrix with (i, j) th entry equal to $B_j''(\tilde{x}_i)$ and \tilde{x}_i is the i th entry of $\tilde{\mathbf{x}} = (\kappa_1, (\kappa_1 + \kappa_2)/2, \kappa_2, \kappa_2, (\kappa_2 + \kappa_3)/2, \kappa_3, \dots, \kappa_{K+7}, (\kappa_{K+7} + \kappa_{K+8})/2, \kappa_{K+8})$. Also, $\boldsymbol{\omega}$ is the $3(K+7) \times 1$ vector given by

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

where $(\Delta\kappa)_k \equiv \kappa_{K+1} - \kappa_k, 1 \leq k \leq K$. The expression in (1.13) involves Simpson's rule, which is applied over each of the inter-knot differences. For more detail on the construction of (1.13) see Section 2 of Wand & Ormerod (2008). Throughout this thesis we work with mixed model representations of Bayesian hierarchical models and thus show how to incorporate O-Sullivan penalised splines using this representation in the following section.

1.10.1 Mixed model representation

The penalised splines can be written in a linear mixed model form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_u^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sigma_\varepsilon^2 \mathbf{I} \end{bmatrix}\right), \quad (1.14)$$

where $\mathbf{X} = [1 \ x_i]_{1 \leq i \leq n}$ and \mathbf{Z} are design matrices and \mathbf{Z} contains B-spline basis functions. A convenient form of expressing $\hat{\boldsymbol{\nu}}$ in (1.12) is by using best linear unbiased prediction (BLUP) of $\boldsymbol{\beta}$ and \mathbf{u} (e.g. Ruppert *et al.*, 2003, Section 4.5.3). This leads to

$$\hat{\boldsymbol{\nu}} = \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{u}} \end{bmatrix} = \left(\mathbf{C}^\top \mathbf{C} + \lambda \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \right)^{-1} \mathbf{C}^\top \mathbf{y}, \quad (1.15)$$

where $\lambda = \sigma_u^2 / \sigma_\varepsilon^2$ and $\mathbf{C} = [\mathbf{X} | \mathbf{Z}]$. The equivalence of (1.15) and (1.12) can be quantified by using a $(K+4) \times (K+4)$ linear transformation matrix \mathbf{L} such that

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

Wand & Ormerod (2008) show that the spectral decomposition of $\boldsymbol{\Omega}$ has the form $\boldsymbol{\Omega} = \mathbf{U} \text{diag}(\mathbf{d}) \mathbf{U}^\top$, where $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ and \mathbf{d} is a $(K+4) \times 1$ vector with 2 zero entries and all others positive. Also, the linear transformation matrix \mathbf{L} has the form $\mathbf{L} = [\mathbf{U}_\mathbf{X} | \mathbf{U}_\mathbf{Z} \text{diag}(\mathbf{d}_\mathbf{Z}^{-1/2})]$, where $\mathbf{d}_\mathbf{Z}$ is the $(K+2) \times 1$ sub-vector of \mathbf{d} containing the positive entries and $\mathbf{U}_\mathbf{Z}$ is the $(K+4) \times (K+2)$ sub-matrix of \mathbf{U} with columns corresponding to the positive entries in \mathbf{d} . Therefore, O-Sullivan penalised splines can be used for fitting (1.14) with the following design matrices:

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

whilst keeping in mind that $\mathbf{B}\mathbf{U}_\mathbf{X}$ is a basis for straight lines (Speed, 1991), and using \mathbf{X} instead will not affect the fit.

1.11 Statistical software

Many semiparametric and nonparametric regression models can be formulated as hierarchical Bayesian models. An upside to this type of formulation is that MCMC software for these type of models can easily be used for fitting and inference. It is possible to tackle different semiparametric regression contexts through the use of BUGS (Bayesian inference Using Gibbs Sampling) inference engine (Spiegelhalter *et al.*, 2003).

In the R computing environment (R Core Team, 2013) BUGS can be accessed through the package `BRugs` (Ligges *et al.*, 2009), thus enabling complete analysis through a single R script. At the moment, `BRugs` is only compatible with the Windows operating system communicating through a version of BUGS known as `OpenBUGS` (Thomas *et al.*, 2006).

MCMC samples can also be obtained using `Stan` (Stan Development Team, 2014) through the R computing environment, made possible using the `rstan` package (Stan Development Team, 2014). `Stan` implements Bayesian inference using Hamiltonian Monte Carlo (HMC) sampling, written in C ++. An alternative to MCMC, which we explore in this thesis, is the variational approximation and in particular MFVB. Some inference engines have come to light recently for performing variational inference using graphical models. These include `VIBES` (Bishop *et al.*, 2003), short for Variational Inference for BayESian networks and `Infer.NET` (Minka *et al.*, 2008). Even though a very useful tool in variational approximations, we do not make use of `VIBES` or `Infer.NET` as they do not cater to the complexity of most of the models considered in this thesis.

In summary, graphical model-based Bayesian inference engines such as `BUGS`, `Stan` and `Infer.NET`, currently entertain a large variety of semiparametric and nonparametric regression settings (e.g. Marley & Wand, 2010; Luts *et al.*, 2014). `BUGS` and `Stan` both make use of Monte Carlo methodology, which is extremely accurate but can be quite slow. `Infer.NET` on the other hand makes use of variational approximations, such as MFVB, which are slightly less accurate but much faster.

1.11.1 Example in `BRugs` and `RStan`

Here we present R code for the simple Binomial model with Beta prior:

$$\begin{aligned} X|p &\sim \text{Binomial}(n, p) \\ p &\sim \text{Beta}(A, B). \end{aligned}$$

After specification of the data, sample size n , probability p , and the two hyperparameters A and B , the model in BRugs takes the form:

```
# Specify model in BRugs:
model
{
  for(i in 1:n)
  {
    x[i] ~ dbern(p)
  }
  p ~ dbeta(A,B)
}
```

Whilst coding the model in RStan is of the form:

```
# Specify model in RStan:
binomBetaModel <- '
  data
  {
    int<lower=1> n;
    int<lower=0,upper=1> x[n];
  }
  parameters
  {
    real<lower=0,upper=1> p;
  }
  model
  {
    for (i in 1:n)
      x[i] ~ bernoulli(p);
    p ~ beta(A,B);binomBetaModel
  }'
```

Throughout this thesis, we make use of both BRugs and RStan for MCMC sampling.

1.12 The flow of the thesis

The primary aim of this thesis is to determine whether MFVB can be used as an alternative to traditional methods for performing inference in semiparametric and nonparametric regression settings. We show that this is possible for various regression settings to allow for fast approximate inference when traditional inference is time consuming or intractable. Each chapter in this thesis is dedicated to the development of an MFVB algorithm catered to a particular model setting. In some instances this requires extension of previously developed MFVB methodology to the model being considered, whereas in others, it involves the original development of an MFVB algorithm. We then use various diagnostic techniques to assess the efficacy of the MFVB algorithm in terms of convergence, timing and accuracy.

Chapter 2 sees the development of an MFVB algorithm for the marginal longitudinal semiparametric regression setting. There are numerous contributions to the topic, however we focus on the Bayesian penalised spline approach of Al Kadiri *et al.* (2010) who use MCMC as an effective way to fitting their class of models and thus we present a comparison of both approaches. The main contribution here lies in the MFVB algorithm being one of the first variational algorithms to involve estimation of an unstructured covariance matrix. In Chapter 3, a modification of MFVB is discussed, known as *non-conjugate variational message passing*, which aids in the incorporation of heteroscedasticity, whilst also involving only closed form algebraic expressions. Chapter 4 looks into the development of MFVB algorithms for three extensions of the heteroscedastic setting given in Chapter 3. This includes an algorithm catered to performing semiparametric regression in real time, bivariate predictor nonparametric regression and in an additive model setting. Accuracy scores and time comparisons are also considered as an effective assessment of the algorithms used. In Chapter 5 we focus on the development of MFVB algorithms for fitting large data sets that exhibit multilevel and longitudinal structures. In particular, we take advantage of a streamlined approach to fitting these type of data. Two types of models are considered: two-level and three-level Gaussian response models. Each are motivated by a corresponding real dataset. Chapter 6 is concerned with MFVB inference for mixture models in measurement error problems. This chapter is heavily motivated by the area of epidemiology, where measurements often do not get measured accurately on all subjects. In Chapter 7, we present a discussion about an alternative approach to MFVB known as

variational message passing, which produces identical results to MFVB but is based on a different algebraic system that allows easier extension to arbitrarily large models. We show this approach for the models considered in Chapter 2 and 6.

Chapter 2

Mean field variational Bayes for marginal longitudinal semiparametric regression

2.1 Introduction

The definitive property of a longitudinal data set is repeated measurements on subjects enabling direct study of change over time. Longitudinal data hence require non-traditional statistical methods, since the set of measurements on one subject tends to be intercorrelated, which must be taken into account when drawing statistical inferences. Over the past decade, a great deal of exposure of longitudinal data in semiparametric regression has emerged. In particular, the special case in which the covariance matrix for each subject's responses is treated as an unspecified parameter to be estimated. This is known as the *marginal longitudinal semiparametric regression* problem. Benefits to marginal longitudinal semiparametric regression analysis include flexible estimation of regression functions and marginal estimation of the covariance matrix of the response vector for each subject.

A proliferation of contributions to the topic emerged in wake of Lin & Carroll (2001) showing that ordinary kernel smoothers are more efficient if *working-independence* is assumed. A summary of research on this problem prior to 2007 is outlined in Ruppert *et al.* (2009, section 3.9). One estimation method for marginal models is MCMC. The primary

The content of this chapter is published as: Menictas, M. and Wand. M.P. (2013). Variational inference for marginal longitudinal semiparametric regression. *Stat*, **2**, 61–71. This research was also presented at the *Young Statisticians Conference*, Melbourne, 2013.

aim of this chapter is to further discussions on the benefit of using a fast deterministic alternative to MCMC known as MFVB, for the marginal longitudinal semiparametric regression problem and some of its semiparametric extensions.

We focus on the Bayesian penalized spline approach of Al Kadiri *et al.* (2010) for function estimation. For fitting, Al Kadiri *et al.* (2010) use MCMC as a fitting and inference tool. However, they admit that such an approach is quite slow in estimation. Their real data example takes almost an hour to run on a contemporary laptop when using BUGS for the MCMC. The variational algorithm developed in this chapter, based on the MFVB paradigm, fits the data in seconds with very similar results. In addition, our algorithm is one of the first variational algorithms involving estimation of an unstructured covariance matrix.

Section 2.2 provides a brief overview of the marginal longitudinal semiparametric regression models considered in Al Kadiri *et al.* (2010). The corresponding MFVB algorithm is developed in Section 2.3. Section 2.4 provides a simulation study for assessing the performance of MFVB against that of MCMC. We provide a real data example in Section 2.5 and the appendices provide algebraic details tied to the methodology in this chapter.

2.2 Marginal longitudinal semiparametric regression models

Nonparametric techniques have recently come into view as an effective way to model longitudinal data. Longitudinal data comprise repeated measurements which are recorded over time on the same subject. We use y_{ij} to denote a measurement for the i th subject at the j th time point. Let $\mathbf{y}_i = (y_{i1}, \dots, y_{in})$ be the vector of responses, $\mathbf{x}_i = (x_{i1}, \dots, x_{in})$ be the vector of predictors for the i th subject, where $1 \leq i \leq m$ and $1 \leq j \leq n$.

As such, the marginal longitudinal nonparametric regression model is

$$E(y_{ij}) = f(x_{ij}), \quad \text{Cov}\{\mathbf{y}_i | f(\mathbf{x}_i)\} = \mathbf{\Sigma}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n, \quad (2.1)$$

where we assume that the mean function f is a real-valued smooth function, and the covariance matrix $\mathbf{\Sigma}$ has dimension $n \times n$. In order to allow for cases where $f(\mathbf{x}_i)$ is random as specified by the model, rather than using $\text{Cov}(\mathbf{y}_i)$ to denote the covariance of \mathbf{y}_i , we use $\text{Cov}\{\mathbf{y}_i | f(\mathbf{x}_i)\}$. Figure 2.1 shows, by example, a simulated data set for model

2.2. MARGINAL LONGITUDINAL SEMIPARAMETRIC REGRESSION MODELS

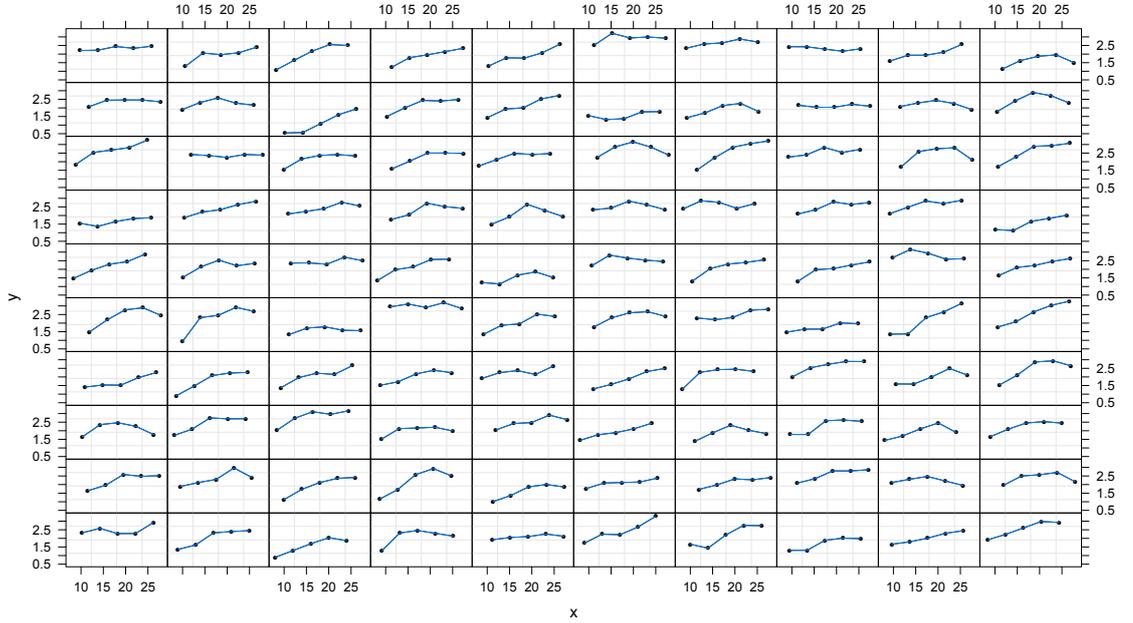


Figure 2.1: A data set simulated from the marginal longitudinal nonparametric regression model (2.1) with $m = 100$, $n = 5$ and f and Σ are as described in (2.2).

(2.1), with $m = 100$, $n = 5$,

$$f(x) = 1.5 + \sin(0.2(x/2 - 4)) \text{ and } \Sigma = \begin{bmatrix} 0.160 & 0.120 & 0.090 & 0.068 & 0.051 \\ 0.120 & 0.160 & 0.120 & 0.090 & 0.068 \\ 0.090 & 0.120 & 0.160 & 0.120 & 0.090 \\ 0.068 & 0.090 & 0.120 & 0.160 & 0.120 \\ 0.051 & 0.068 & 0.090 & 0.120 & 0.160 \end{bmatrix}. \quad (2.2)$$

We are interested in efficient estimation of f and Σ from data such as that shown in Figure 2.1. The procedure we use to obtain the mean function estimate involves spline models for f of the form

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k z_k(x), \quad (2.3)$$

where z_k , $1 \leq k \leq K$ represents a set of suitably transformed cubic O'Sullivan splines, as described in Section 1.10. The quantity of spline basis functions K has a small effect on the adequacy of (2.3). However, in practice, $K = 25$ is a sufficient choice for most spline basis functions (Li & Ruppert, 2008). In order to prevent over fitting of the function, we penalize the spline basis function coefficients u_k , $1 \leq k \leq K$, by treating them as a random sample from a distribution which has a mean of zero and a variance of σ^2 . This allows for

the following linear mixed model representation of (2.1) and (2.3):

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

where

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_1 \\ \vdots & \vdots \\ \mathbf{1} & \mathbf{x}_m \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} z_1(\mathbf{x}_1) & \dots & z_K(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ z_1(\mathbf{x}_m) & \dots & z_K(\mathbf{x}_m) \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \vdots \\ \boldsymbol{\varepsilon}_m \end{bmatrix},$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \quad \text{and} \quad \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_K \end{bmatrix}.$$

The random effects have a mean of zero and a covariance matrix:

$$\text{Cov} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon}_1 \\ \boldsymbol{\varepsilon}_2 \\ \vdots \\ \boldsymbol{\varepsilon}_m \end{bmatrix} = \begin{bmatrix} \sigma^2 \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \boldsymbol{\Sigma} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Sigma} \end{bmatrix} = \begin{bmatrix} \sigma^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma} \end{bmatrix}.$$

The model for the distribution of \mathbf{u} and $\boldsymbol{\varepsilon}$ is

$$\begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma} \end{bmatrix} \right). \quad (2.5)$$

Sections 2.2.1, 2.2.2 and 2.2.3 provide semiparametric extensions to model (2.1).

2.2.1 Additive models extension

There are various semiparametric regressions extensions of (2.1) that could be considered. Here we consider the additive model extension. This involves the incorporation of several continuous predictor variables for each scalar response variable. We restrict our discussion to the incorporation of two continuous predictor variables x_{1ij} and x_{2ij} corresponding to each y_{ij} . In this case, the marginal longitudinal additive model is

$$\begin{aligned} E(y_{ij}) &= \beta_0 + f_1(x_{1ij}) + f_2(x_{2ij}), \\ \text{Cov}\{y_i | f_1(\mathbf{x}_{1i}), f_2(\mathbf{x}_{2i})\} &= \boldsymbol{\Sigma}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n, \end{aligned} \quad (2.6)$$

where we assume that the mean functions f_1 and f_2 are smooth functions and the covariance matrix $\boldsymbol{\Sigma}$ has dimension $n \times n$. The procedure we use to obtain the mean function

estimates involve spline models for f_1 and f_2 of the following form:

$$f_1(x_1) = \beta_{11}\mathbf{x}_1 + \sum_{k=1}^{K_1} u_{1k}z_{1k}(\mathbf{x}_1) \quad \text{and} \quad f_2(\mathbf{x}_2) = \beta_{21}x_2 + \sum_{k=1}^{K_2} u_{2k}z_{2k}(\mathbf{x}_2), \quad (2.7)$$

with coefficients $u_{1k} \stackrel{\text{ind.}}{\sim} N(0, \sigma_1^2)$, $1 \leq k \leq K_1$ and $u_{2k} \stackrel{\text{ind.}}{\sim} N(0, \sigma_2^2)$, $1 \leq k \leq K_2$. This allows for the following Gaussian linear mixed model representation of (2.6) and (2.7):

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}.$$

The differences between this model and the one considered in Section 2.2 lie in the structure of the design matrices. The design matrices are now

$$\mathbf{X} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{11} & \mathbf{x}_{21} \\ \vdots & \vdots & \vdots \\ \mathbf{1} & \mathbf{x}_{1m} & \mathbf{x}_{2m} \end{bmatrix}, \quad \text{and}$$

$$\mathbf{Z} = \begin{bmatrix} z_{11}(\mathbf{x}_{11}) & \dots & z_{1K_1}(\mathbf{x}_{11}) & z_{21}(\mathbf{x}_{21}) & \dots & z_{2K_2}(\mathbf{x}_{21}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ z_{11}(\mathbf{x}_{1m}) & \dots & z_{1K_1}(\mathbf{x}_{1m}) & z_{21}(\mathbf{x}_{2m}) & \dots & z_{2K_2}(\mathbf{x}_{2m}) \end{bmatrix},$$

where \mathbf{x}_{1i} is the $n \times 1$ vector containing the x_{1ij} measurements and \mathbf{x}_{2i} is the $n \times 1$ vector containing the x_{2ij} measurements. The coefficient vectors are

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_{11} \\ \beta_{12} \end{bmatrix}, \quad \text{and} \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix},$$

where \mathbf{u}_1 is the $K_1 \times 1$ vector containing the u_{1k} and \mathbf{u}_2 is the $K_2 \times 1$ vector containing the u_{2k} . A convenient assumption for estimating σ_1^2 , σ_2^2 , and $\boldsymbol{\Sigma}$ is now

$$\begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \boldsymbol{\varepsilon} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_1^2 \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \sigma_2^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I}_m \otimes \boldsymbol{\Sigma} \end{bmatrix} \right). \quad (2.8)$$

Fitting via MFVB is analogous to that described in Section 2.3. The difference here is that there are two variance parameters σ_1^2 and σ_2^2 (where in extensions to additive models with d smooth functions there will be d such variance parameters) as well as the error covariance matrix $\boldsymbol{\Sigma}$. A simpler type of additive model is

$$E(y_{ij}) = \beta_0 + \beta_1 x_{1ij} + f_2(x_{2ij}), \quad \text{Cov}\{y_i | \mathbf{x}_{1i}, f_2(\mathbf{x}_{2i})\} = \boldsymbol{\Sigma}, \quad (2.9)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n.$$

This model represents a semiparametric regression model in the true sense, since the right-hand side comprises a parametric component $\beta_1 x_{1ij}$ where the effect of the x_{1ij} s are modelled parametrically *and* a nonparametric component $f_2(\mathbf{x}_{2i})$, where the effect of the x_{2ij} s are modelled nonparametrically. This is a simpler model to that of (2.6), since there is only one smooth function component. However, the linear mixed model corresponding to this model is similar to that in Section 2.2, where the random component structure (2.5) applies to (2.9).

2.2.2 Incorporation of interactions

Given the situation where an additive model is not guaranteed to provide a satisfactory fit of the data, one may opt to incorporate an interaction term. Such an extension of (2.6) is

$$\begin{aligned} E(y_{ij}) &= \beta_0 + f_{x_{3ij}}(x_{1ij}) + f_2(x_{2ij}), \\ \text{Cov}\{\mathbf{y}_i | f_{x_{3ij}}(\mathbf{x}_{1i}), f_2(\mathbf{x}_{2i})\} &= \boldsymbol{\Sigma}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n, \end{aligned} \tag{2.10}$$

where the x_{3ij} correspond to measurements on the categorical variable x_3 . For example, if $x_{3ij} \in \{0, 1\}$ then we have $E(y_{ij}) = \beta_0 + f_0(x_{1ij}) + f_2(x_{2ij})$ if $x_{3ij} = 0$ and $E(y_{ij}) = \beta_1 + f_1(x_{1ij}) + f_2(x_{2ij})$ if $x_{3ij} = 1$. Interaction models such as (2.10) have two important uses. One of which is the ability to check for additivity in the model. This can be done by testing the null hypothesis of a true additive model i.e., one without any interaction terms, against the alternative hypothesis of the additive model involving interaction terms. If the null hypothesis is accepted, we assume that using the additive model will provide a reasonable fit to the data. Another important use of interaction models is the ability to use alternative models when an additive model fails to fit the data well. When the null hypothesis of an additive model is rejected, appropriate interaction terms are added to the additive model.

The Gaussian linear mixed model for fitting (2.10) has the same structure as that of fitting the additive model (2.6), where the only difference is the form of the design matrices \mathbf{X} and \mathbf{Z} . An in-depth explanation including examples may be found in Ruppert *et al.* (2003).

2.2.3 Varying coefficient models

Varying coefficient models are another type of multiple-predictor semiparametric regression model, where the coefficients are allowed to vary as smooth functions of other variables. The model

$$E(y_{ij}) = f_0(s_{ij}) + f_1(s_{ij})x_{ij}, \quad \text{Cov}\{\mathbf{y}_i | f_0(\mathbf{s}_i)\} = \boldsymbol{\Sigma}, \quad (2.11)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n,$$

is the simplest marginal longitudinal varying coefficient model, where s_{ij} are the longitudinal measurements on a continuous predictor variable s and x_{ij} are the measurements on a second predictor x . Model (2.11) shows that s has a modifying effect on the linear relationship between $E(\mathbf{y})$ and x , which is modelled through the varying coefficients $f_0(s)$ and $f_1(s)$. The gaussian linear mixed model for fitting (2.11), has the same structure as that for fitting the additive model (2.6), where the only difference is the form of the design matrices \mathbf{X} and \mathbf{Z} . Relevant details are provided in Section 12.4 of Ruppert *et al.* (2003).

2.3 Mean field variational Bayes methodology

Each of the marginal longitudinal semiparametric regression models in section 2.2, and their extensions to d smooth functions, can be treated using the Gaussian linear mixed model

$$\mathbf{y} | \mathbf{u}, \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}), \quad \mathbf{u} | \boldsymbol{\sigma}^2 \sim N\left(\mathbf{0}, \text{blockdiag}\left(\sigma_\ell^2 \mathbf{I}_{K_\ell}\right)_{1 \leq \ell \leq d}\right), \quad (2.12)$$

where K_ℓ corresponds to the number of spline basis functions used in the ℓ th smooth function estimate. Maximum likelihood and best prediction can be used as inference strategies for constructing estimates of the regression function f . Al Kadiri *et al.* (2010) discuss that, despite the Gaussian linear mixed model (2.4) being a simple model, fitting via standard mixed model software such as `lme()` (Pinheiro *et al.*, 2009) in the R computing language has not yet been successful. This led Al Kadiri *et al.* (2010) to consider the Bayesian inference version and implementation via MCMC. In this section, we discuss and define a deterministic approximation alternative to MCMC known as mean field variational Bayes. MFVB scales well to large applications with simple implementation and requires less computation time than MCMC. However, MFVB cannot generate exact results and seems to be limited in its approximation accuracy. Nevertheless, in certain circumstances

MFVB may be *very* accurate.

We consider working with a hierarchical Bayesian version of (2.12), where $\boldsymbol{\beta}(p \times 1)$, $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_d^2)$ and $\boldsymbol{\Sigma}(n \times n)$ are treated as random. Common prior distributions for these parameters are:

$$\boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{F}), \quad \sigma_\ell \sim \text{Half-Cauchy}(A_\ell), \quad \boldsymbol{\Sigma} \sim \text{Inverse-Wishart}(A_\Sigma, B_\Sigma). \quad (2.13)$$

An important concern to fitting and inference is the choice of the hyperparameters \mathbf{F} , A_ℓ , A_Σ and B_Σ . When prior beliefs about the model parameters are held, quantities may be used to represent these beliefs about the hyperparameters. However in most instances, prior beliefs are unavailable. In cases such as these one may use vague priors. We assume that the data have been standardised before we impose the following priors for both the fixed effects and variance hyperparameters:

$$\mathbf{F} = \mathbf{I} \times 10^{10}, \quad A_\ell = 1 \times 10^5, \quad A_\Sigma = n, \quad \text{and} \quad B_\Sigma = 0.01\mathbf{I}_n. \quad (2.14)$$

Approximate inference via MCMC is aided by the distribution theoretical result given in Result 1.4.6. We achieve the following equivalent form of (2.12) and (2.13):

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma} &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}), \quad \mathbf{u} | \sigma^2 \sim N\left(\mathbf{0}, \text{blockdiag}_{1 \leq \ell \leq d}(\sigma_\ell^2 \mathbf{I}_{K_\ell})\right), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \mathbf{F}), \quad \sigma_\ell^2 | a_\ell \sim \text{Inverse-Gamma}\left(\frac{1}{2}, \frac{1}{a_\ell}\right), \quad a_\ell \sim \text{Inverse-Gamma}\left(\frac{1}{2}, \frac{1}{A_\ell^2}\right), \\ \boldsymbol{\Sigma} &\sim \text{Inverse-Wishart}(A_\Sigma, B_\Sigma), \end{aligned} \quad (2.15)$$

where A_ℓ and A_Σ are positive constants, \mathbf{F} and B_Σ are both positive definite matrices. Figure 2.2(a) illustrates the DAG corresponding to (2.15) and is effective in determining its full conditional distributions. The posterior distributions of $\boldsymbol{\beta}$, \mathbf{u} , σ_ℓ , $\boldsymbol{\Sigma}$, and a_ℓ , are not available in closed form. However the distributions of the full conditionals are shown to be:

$$\begin{aligned} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big|_{\text{rest}} &\sim N\left(\left(\mathbf{C}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma})^{-1} \mathbf{C} + \begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq d} \left(\frac{1}{\sigma_\ell^2} \mathbf{I}_{K_\ell}\right) \end{bmatrix}\right)^{-1} \mathbf{C}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}) \mathbf{y}, \right. \\ &\quad \left. \left(\mathbf{C}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}) \mathbf{C} + \begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq d} \left(\frac{1}{\sigma_\ell^2} \mathbf{I}_{K_\ell}\right) \end{bmatrix}\right)^{-1}\right), \end{aligned}$$

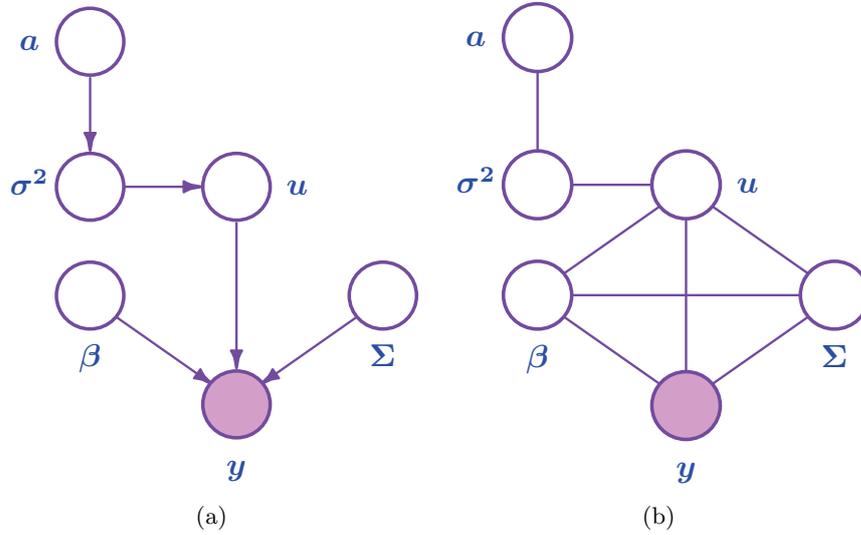


Figure 2.2: (a) Directed acyclic graph representation of the Bayesian penalized spline model (2.15), where $\mathbf{a} = (a_1, \dots, a_d)$, and $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_d^2)$. (b) Moral graph representation of the DAG in (a).

$$\sigma_\ell^2 | \text{rest} \sim \text{Inverse-Gamma} \left(\frac{1}{2} (K_\ell + 1), \frac{1}{2} \|u_\ell\|^2 + a_\ell^{-1} \right),$$

$$a_\ell | \text{rest} \sim \text{Inverse-Gamma} \left(1, (\sigma_\ell^{-2} + A^{-2}) \right) \text{ and}$$

$$\Sigma | \text{rest} \sim \text{Inverse-Wishart} \left(A_\Sigma + m, \sum_{i=1}^m (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{u}) (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{u})^\top + B_\Sigma \right).$$

An elementary form of MCMC sampling, called Gibbs sampling, can be used to draw successive samples from these full conditional distributions. However, we demonstrate fitting via MFVB as it provides for faster fitting and inference to MCMC.

We next consider implementation of MFVB methodology, which involves restricting the full posterior density function to have the approximate product form

$$p(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}, \boldsymbol{\sigma}^2, \Sigma | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}) q(\boldsymbol{\sigma}^2, \Sigma). \quad (2.16)$$

Additional factorizations, referred to as *induced factorizations*, arise from an interaction between the factorization assumed in the variational posterior distribution and the conditional independence properties of the true joint distribution (Bishop, 2006). Such induced factorizations can easily be detected using moralisation (described in Sections 1.8.2 and 1.8.3). We show this using the moral graph of Figure 2.2(b). Let $A = \sigma^2$, $B = \Sigma$, 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 from the moral graph in Figure 2.2(b) that all paths between σ^2 and Σ must

pass through at least one of $\{\boldsymbol{\beta}, \mathbf{u}, \mathbf{y}\}$. Thus, applying Theorem 1.8.1 we see that

$$\sigma^2 \perp\!\!\!\perp \boldsymbol{\Sigma} \mid \{\boldsymbol{\beta}, \mathbf{u}, \mathbf{y}\}.$$

Similarly, we can say that $\{\boldsymbol{\beta}, \mathbf{u}\} \perp\!\!\!\perp \mathbf{a} \mid \{\mathbf{y}, \sigma^2, \boldsymbol{\Sigma}\}$, and hence (2.16) is updated to

$$q(\boldsymbol{\beta}, \mathbf{u}, \mathbf{a}, \sigma^2, \boldsymbol{\Sigma}) = q(\boldsymbol{\beta}, \mathbf{u})q(\mathbf{a})q(\sigma^2)q(\boldsymbol{\Sigma}).$$

Then, as shown through the derivations in Appendix 2.A,

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}) &\text{ is a Multivariate Normal density function,} \\ q^*(\sigma_\ell^2) &\text{ and } q^*(a_\ell) \text{ are each Inverse Gamma density functions,} \\ q^*(\boldsymbol{\Sigma}) &\text{ is an Inverse-Wishart density function.} \end{aligned}$$

Note that $\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ denote the mean vector and covariance matrix for $q^*(\boldsymbol{\beta}, \mathbf{u})$, and $A_{q(\sigma_\ell^2)}$ and $B_{q(\sigma_\ell^2)}$ denote the shape and rate parameters for $q^*(\sigma_\ell^2)$. Similar definitions apply for the other density functions. The optimal values of these parameters are determined using Algorithm 1. The lower bound on the marginal log-likelihood is (the

Set up initial values:

$$\mu_{q(1/\sigma_\ell^2)} > 0, \mu_{q(1/a_\ell)} > 0, 1 \leq \ell \leq d, M_{q(\boldsymbol{\Sigma}^{-1})} \text{ positive definite.}$$

Cycle through:

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \left\{ \mathbf{C}^\top \left(\mathbf{I}_m \otimes \mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right) \mathbf{C} + \begin{pmatrix} \mathbf{F}^{-1} & 0 \\ 0 & \text{blockdiag}_{1 \leq \ell \leq d} \left(\mu_{q(1/\sigma_\ell^2)} \mathbf{I}_{K_\ell} \right) \end{pmatrix} \right\}^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \left(\mathbf{I}_m \otimes \mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right) \mathbf{y}$$

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

$$B_{q(a_\ell)} \leftarrow \mu_{q(1/\sigma_\ell^2)} + A^{-2} \quad ; \quad \mu_{q(1/a_\ell)} \leftarrow 1/B_{q(a_\ell)}$$

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

$$\mu_{q(1/\sigma_\ell^2)} \leftarrow \frac{1}{2}(K_\ell + 1)/B_{q(\sigma_\ell^2)}$$

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

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

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

Algorithm 1: *Mean field variational Bayes algorithm for determination of the optimal parameters in $q^*(\boldsymbol{\beta}, \mathbf{u})$, $q^*(\sigma^2)$, $q^*(\mathbf{a})$, and $q^*(\boldsymbol{\Sigma})$.*

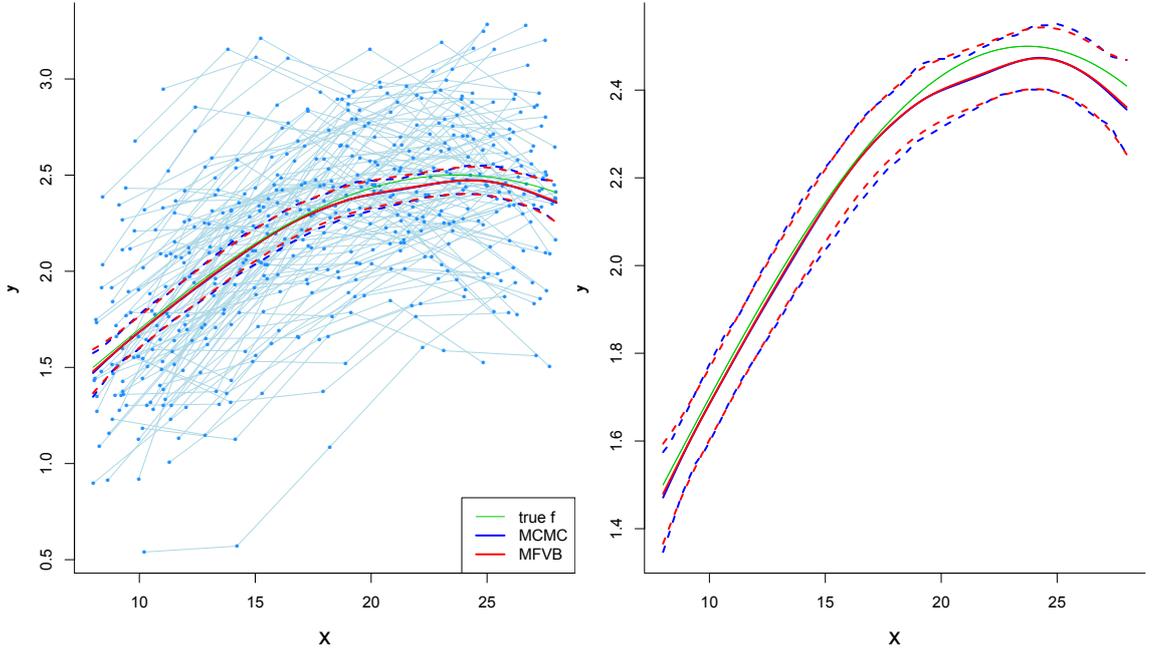


Figure 2.3: *Fitted curve estimates and pointwise 95% credible sets for both MCMC and MFVB approaches for fitting the Bayesian penalized spline model (2.15) to simulated data. Right panel: zoomed view with data points removed.*

derivation of which can be found in Appendix 2.B):

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} \left(\sum_{\ell=1}^d K_{\ell} + p \right) - \frac{1}{2} m \log(2\pi) - \log(\pi) + \sum_{\ell=1}^d \log \Gamma\left(\frac{1}{2}(K_{\ell} + 1)\right) \\
 &\quad - \frac{1}{2} \log |\mathbf{F}| - \log(A) - \frac{1}{2} \text{tr} \left\{ \mathbf{F}^{-1} \left(\boldsymbol{\mu}_{q(\beta)} \boldsymbol{\mu}_{q(\beta)}^{\top} + \boldsymbol{\Sigma}_{q(\beta)} \right) \right\} \\
 &\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, u)}| - \frac{1}{2} \sum_{\ell=1}^d (K_{\ell} + 1) \log(B_q(\sigma_{\ell}^2)) \\
 &\quad - \frac{1}{2} (A_{\boldsymbol{\Sigma}} + m) \log |B_q(\boldsymbol{\Sigma})| - \sum_{\ell=1}^d \log \left(\mu_{q(1/\sigma_{\ell}^2)} + A^{-2} \right) \\
 &\quad + \sum_{\ell=1}^d \mu_{q(1/\sigma_{\ell}^2)} \mu_{q(1/a_{\ell})} + \frac{1}{2} A_{\boldsymbol{\Sigma}} (\log |B_{\ell}|) - \log(C_{n, A}) \\
 &\quad + \log(C_{n, A+m}).
 \end{aligned} \tag{2.17}$$

Figure 2.3 illustrates the fitted curves using both MCMC and MFVB for the Bayesian penalized spline model (2.15). The data were simulated from a version of the marginal longitudinal nonparametric regression model (2.1) with $m = 100$, $n = 5$ and f and $\boldsymbol{\Sigma}$ as described in (2.2). The x_{ij} are equally spaced but with the starting positions x_{i1} generated uniformly from the interval $(8, 8 + \frac{20}{n})$ and $\boldsymbol{\varepsilon} \sim N(0, 1)$. We used diffuse priors given by (2.14). For the MCMC approach, samples of size 10,000 were generated where the first 5,000 values were discarded. For the MFVB approach the iterations were terminated when the relative change in $\log \underline{p}(\mathbf{y}; q)$ fell below 10^{-7} . For this example the MCMC and MFVB fits and pointwise 95% credible sets are almost identical. This suggests that MFVB

achieves high accuracy for Gaussian response Bayesian penalized spline regression.

2.3.1 Heuristic justification of mean field variational Bayes

Here we are interested in justifying the initial MFVB approximate product restriction given in (2.16) for the model in (2.15). We consider the marginal longitudinal frequentist regression model

$$\mathbf{y} \sim \text{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}). \quad (2.18)$$

The likelihood-based estimates and confidence intervals attained from (2.18) will be similar to the Bayes estimates and credible sets based on the marginal longitudinal Bayesian regression model

$$\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\Sigma} \sim \text{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}),$$

given that the priors on $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ are noninformative. In addition, asymptotic independence (or parameter orthogonality) of $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ in (2.18) is implied by the block-diagonal form of the Fisher information matrix:

$$\mathbf{I}(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \begin{bmatrix} \mathbf{X}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{m}{2} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \end{bmatrix}. \quad (2.19)$$

The derivation of (2.19) is given in Appendix 2.C. In the *Bayesian world*, such asymptotic independence corresponds to the approximate product form

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma}|\mathbf{y}) \approx p(\boldsymbol{\beta}|\mathbf{y})p(\boldsymbol{\Sigma}|\mathbf{y}),$$

where the MFVB approximate product restriction $q(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = q(\boldsymbol{\beta})q(\boldsymbol{\Sigma})$ corresponds to the replacement of approximate with exact equality in this expression. Approaches such as best prediction (Robinson, 1991) need to be considered for frequentist inference of (2.15) for the estimation of the random effects vector \mathbf{u} . Even so, orthogonality between the mean and variance parameters of the model is still maintained. Hence, we would expect the MFVB approximate product restriction in (2.16) to be reasonable and not incur heavy losses in accuracy.

2.3.2 Mean field variational Bayes approximate density functions for entries of Σ

In Section 2.4, we present the results from a comprehensive simulation study that compares the MFVB methodology against MCMC. We do this comparison on the mean function hexiles, $f(H_1), \dots, f(H_5)$ and the entries of the variance covariance matrix, $\Sigma_{11}, \Sigma_{12}, \dots, \Sigma_{45}$. The mean function hexiles are straightforward to attain, however, the entries of the variance covariance matrix require additional effort. The derivations of the distributions of these entries are given in this section.

2.3.2.1 Diagonal entries

The distributions of the diagonal entries of $q^*(\Sigma)$ are obtained using Theorem 2.3.1.

Theorem 2.3.1. *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:*

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

Details of the proof of Theorem 2.3.1 are given in Appendix 2.D. From Algorithm 1 we see that the optimal q -density of Σ has the form

$$q^*(\Sigma) \text{ has an Inverse-Wishart } (A_{q(\Sigma)}, \mathbf{B}_{q(\Sigma)}) \text{ density function,}$$

where

$$\begin{aligned} A_{q(\Sigma)} &= a + m, \text{ and} \\ \mathbf{B}_{q(\Sigma)} &= \mathbf{B} + (\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\beta, \mathbf{u})})(\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\beta, \mathbf{u})})^\top + \mathbf{C}\Sigma_{q(\beta, \mathbf{u})}\mathbf{C}^\top. \end{aligned}$$

Now, using Theorem 2.3.1, the diagonal entries take the form

$$q^*(\Sigma_{jj}) \sim \text{Inverse-Gamma}\left(\frac{1}{2}(a + m - n + 1), \frac{1}{2}\left\{\mathbf{B} + (\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\beta, \mathbf{u})})(\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\beta, \mathbf{u})})^\top + \mathbf{C}\Sigma_{q(\beta, \mathbf{u})}\mathbf{C}^\top\right\}_{jj}\right),$$

where $j \in 1, \dots, n$.

2.3.2.2 Off-diagonal entries

We considered two main approaches to finding the distributions of the off-diagonal entries of $q^*(\Sigma)$. These included: (*Approach 1*) Numerical characteristic function inversion based on the Fast Fourier Transform; or (*Approach 2*) Monte Carlo sampling of the off-diagonal elements in $q^*(\Sigma)$. It is most desirable to attain these densities exactly, however when exploring *Approach 1* we found that there was no straightforward way of obtaining the characteristic function of an Inverse-Wishart random matrix. As a result, we used *Approach 2*, subject to Monte Carlo error, to finding the distributions of these off-diagonal elements.

Whilst initially exploring *Approach 1* however, we came across some exciting results involving the characteristic function of the off-diagonal elements of a Wishart random matrix. This is presented in Theorem 2.3.2.

Theorem 2.3.2. *The characteristic function of the off-diagonal entries of a Wishart(a, \mathbf{B}) random matrix \mathbf{X} take the form:*

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

where the λ_{ℓ} 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.

The proof of Theorem 2.3.2 is given in Appendix 2.E. This Theorem enables one to find the distributions of the off-diagonal entries of a Wishart random matrix analytically by making use of the Fast Fourier Transform inversion of the characteristic function.

2.4 Simulation study

A simulation study was carried out for Algorithm 1 in order to assess the accuracy of its Bayesian inference compared to that of MCMC. 1000 data-sets were generated, which corresponds to the simulation setting described in Section 5.1 of Al Kadiri *et al.* (2010). The nonparametric regression model is of the form

$$E(y_{ij}|\mathbf{u}) = f(x_{ij}), \quad \text{Cov}(\mathbf{y}_i|\mathbf{u}) = \Sigma, \quad 1 \leq i \leq 100, \quad 1 \leq j \leq 5, \quad (2.20)$$

with

$$f(x_i) = 1 + \frac{1}{2}\Phi\left(\frac{2x - 36}{5}\right) \quad \text{and} \quad \Sigma = \begin{bmatrix} 0.122 & 0.098 & 0.078 & 0.063 & 0.050 \\ 0.098 & 0.122 & 0.098 & 0.078 & 0.063 \\ 0.078 & 0.098 & 0.122 & 0.098 & 0.078 \\ 0.063 & 0.078 & 0.098 & 0.122 & 0.098 \\ 0.050 & 0.063 & 0.078 & 0.098 & 0.122 \end{bmatrix},$$

where Φ is the standard normal distribution function and (2.20) is an example of the special case of the Bayesian hierarchical model (2.15) with $d = 1$.

There are various approaches by which the accuracy of a variational Bayes approximate density function $q^*(\theta)$ may be assessed with respect to the exact posterior density $p(\theta|\mathbf{y})$. Avoiding approximate inference methods would be ideal, however, using MCMC with sufficiently large samples can easily be used to approximate $p(\theta|\mathbf{y})$ quite well. A detailed explanation is provided in Section 1.9.

Many parameters are of interest when conducting an assessment of MFVB against MCMC. We have chosen to work with the mean function hexiles, $f(H_1), \dots, f(H_5)$ and the entries of the variance covariance matrix, $\Sigma_{11}, \Sigma_{12}, \dots, \Sigma_{45}$. The mean function hexiles are straightforward to attain, however, the entries of the variance covariance matrix require additional effort. The derivations of the distributions of these entries were given in the previous section. The accuracy assessments that follow are based on MCMC samples using BUGS iterations of size 5,000. A thinning factor of 5 was applied to post burn-in samples of size 5,000. This resulted in MCMC samples of size 1,000 for density estimation. The density estimates were obtained using the binned kernel density estimate `bkde()` function in the R package `Kernsmooth` (Wand & Ripley, 2009). The MFVB iterations were terminated when the increase in $\log p(\mathbf{y}; q)$ fell below 10^{-7} . Figure 2.4 summarizes the accuracy results whilst Figure 2.5 plots the variational Bayes and MCMC approximate posteriors for the first simulated dataset. The mean function at each hexile and each entry of Σ shows high accuracy, with most accuracy levels above 85%.

Another crucial type of accuracy assessment is the comparison between the advertised coverage of variational Bayes approximate credible intervals and the true coverage. Table 2.1 shows the percentages of the true parameter coverage for the approximate 95% credible intervals formed from the variational Bayes posterior densities with 0.025 probability mass in each tail. The coverage overall is very good and does not fall below 93%.

2.4. SIMULATION STUDY

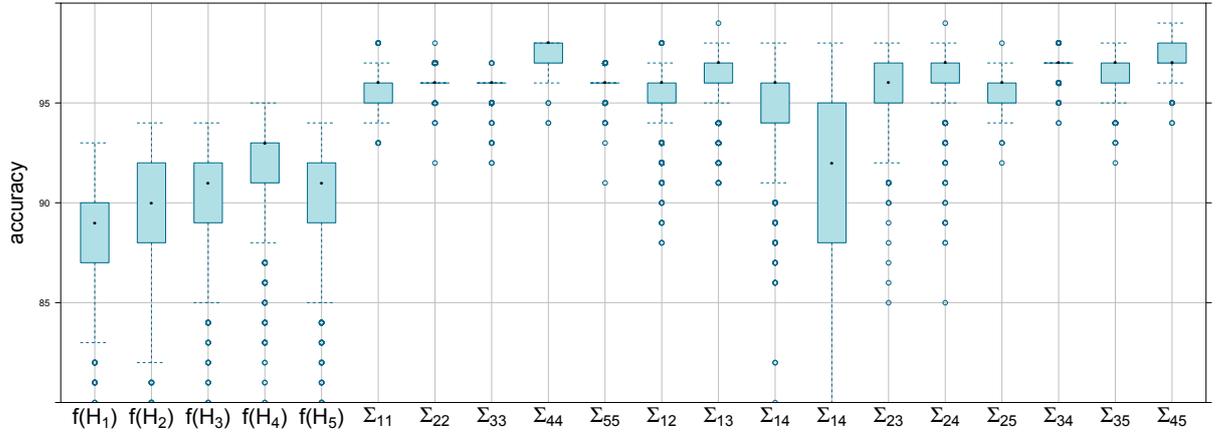


Figure 2.4: Summary of the simulation for the marginal longitudinal nonparametric regression model (2.1), where $f(H_i)$ is the mean function at hexile i . For the mean function at each hexile and each entry of Σ , the accuracy values are summarized as a boxplot.

Σ_{11}	Σ_{22}	Σ_{33}	Σ_{44}	Σ_{55}	Σ_{12}	Σ_{13}	Σ_{14}	Σ_{15}	Σ_{23}
95%	95%	96%	94%	95%	95%	94%	94%	95%	95%
Σ_{24}	Σ_{25}	Σ_{34}	Σ_{35}	Σ_{45}	$f(H_1)$	$f(H_2)$	$f(H_3)$	$f(H_4)$	$f(H_5)$
95%	95%	94%	95%	94%	95%	95%	95%	93%	94%

Table 2.1: Percentage coverage of true parameter values by approximate 95% credible intervals based on variational Bayes approximate posterior density functions. The percentages are based on 1000 replications.

2.4.1 Assessment of speed

To assess the gain in time savings by using MFVB, we monitored the time taken to fit each model based on a different replication in the simulation study. The computations for this study were performed on a laptop computer (Mac OS X; 2.8 GHz processor, 16 GBytes random access memory). The computation times for each approach are summarised in Table 2.2.

MCMC	MFVB
282.555 (4.506)	3.064 (0.278)

Table 2.2: Average (standard deviation) times in seconds for the simulation study comparing MCMC and MFVB fitting of the model in (2.20).

The average computing time for MFVB is about 3 seconds which is substantially faster

2.4. SIMULATION STUDY

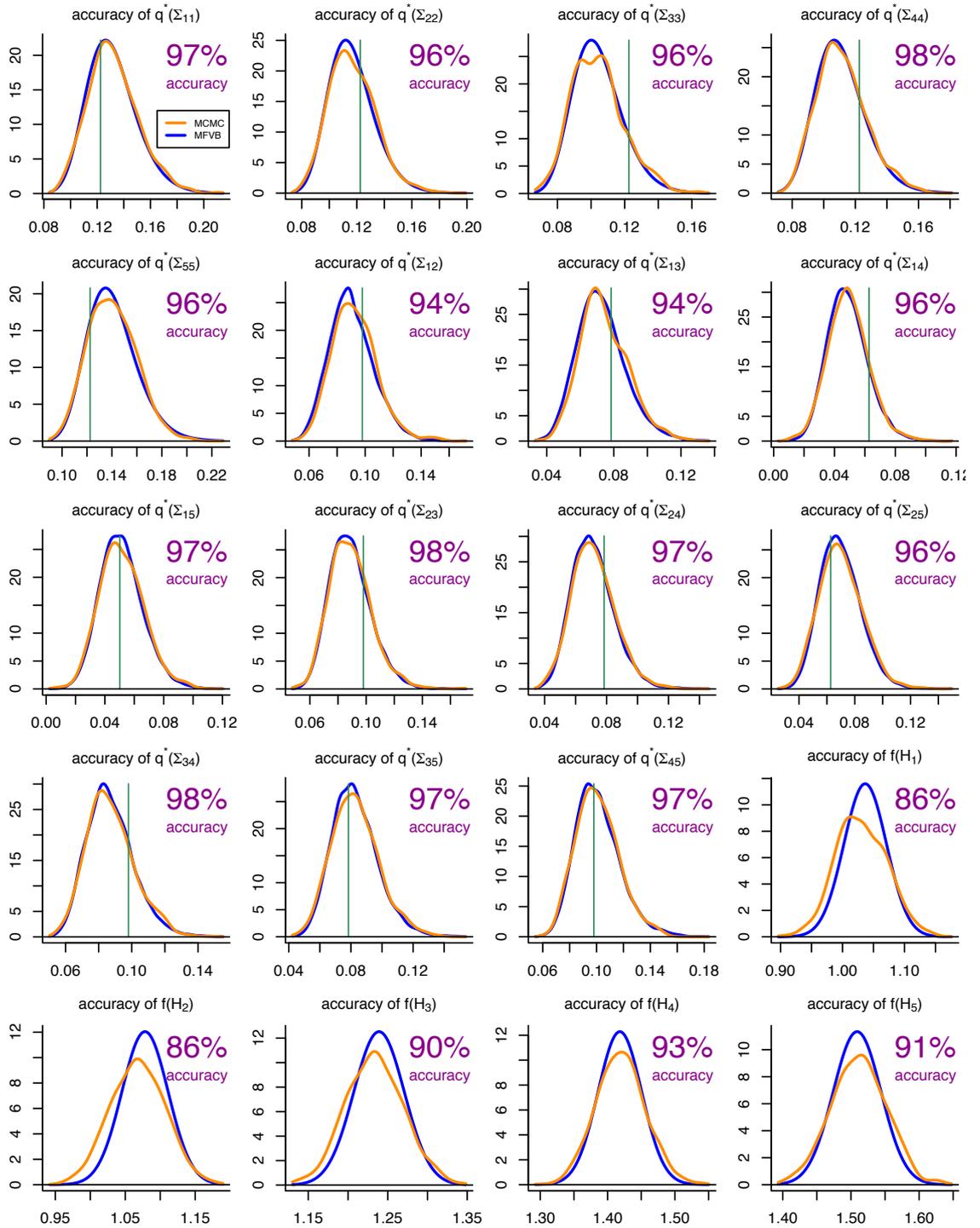


Figure 2.5: Variational Bayes and MCMC approximate posterior densities for the first simulated data set.

than the 4.7 minutes taken by MCMC i.e., 94 times faster. It is evident that the speed gains achieved by MFVB needs to be traded off against the accuracy losses which are incurred by restriction (2.16).

Making time comparisons completely fair between MFVB and MCMC is difficult, since convergence for each approximation method is different. The number of iterations used for MFVB and the number of samples used for MCMC was sufficient for convergence in this particular simulation study setting. However, these values will change depending on the type of application.

2.5 Application

The Bayesian additive/interaction model (2.10) was fitted to data from a nutritional epidemiology study (source: Kipnis et al., 2003), which was presented in Section 5.2 of Al Kadiri *et al.* (2010). Convergence for both MCMC and MFVB was assessed in the same way as given in the previous section. Keeping with the notation of (2.10), the variables are:

- y = logarithm of intake of protein as measured by the biomarker urinary nitrogen,
- x_1 = body mass index,
- x_2 = logarithm of intake of protein as measured by a 24-hour recall instrument,
- x_3 = gender.

Figure 2.6 shows the MCMC and MFVB estimates of f_{male} , f_{female} and f_2 and the corresponding 95% pointwise credible sets. The MCMC fit took about 1 hour and 5 minutes,

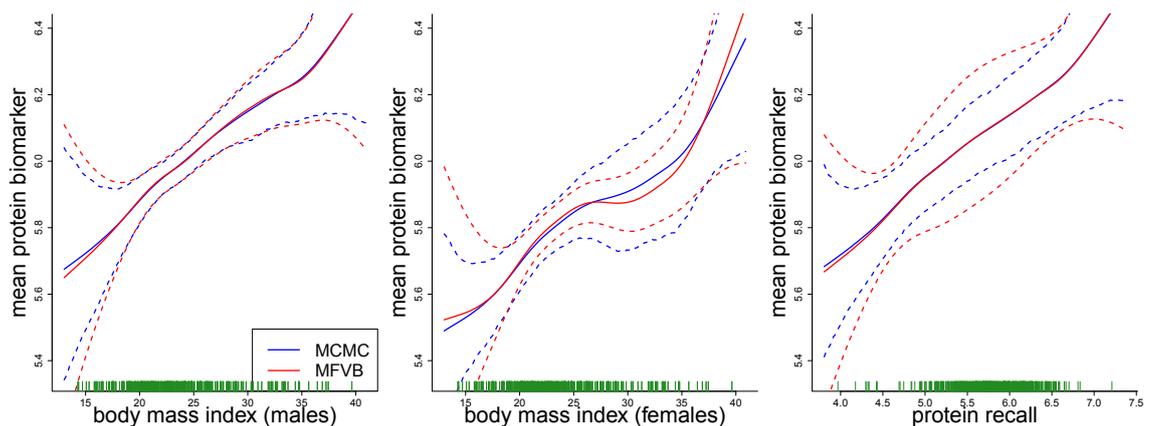


Figure 2.6: *Estimated functions for the additive/interaction model for the nutritional epidemiology data using MCMC and MFVB. The dashed curves correspond to 95% pointwise credible sets.*

whilst the MFVB fit took 20 seconds.

2.6 Discussion

The aim of this chapter was to develop a fast deterministic approach to MCMC to shorten estimation time. We also sought to address the marginal longitudinal regression problem and some of its semiparametric extensions. As a result we have derived a MFVB algorithm for fast approximate inference in semiparametric regression. The central finding here is that, for using the marginal longitudinal semiparametric regression model, MFVB achieves excellent accuracy when compared to MCMC for the main parameters of interest. In addition, the algorithm is one of the first variational algorithms to involve estimation of an unstructured covariance matrix. The real data example presented in Al Kadiri *et al.* (2010) takes almost an hour for MCMC to run on a contemporary laptop using BUGS. The MFVB algorithm developed in this chapter, fits the same data in seconds with very similar results. The simulation study shows evidence of significant speed and accuracy attributes of MFVB against the MCMC benchmark.

2.A Derivation of algorithm 1

Algorithm 1 depends on the following derivations of the optimal density functions. Constants with respect to the function argument are denoted by ‘const’. The MFVB calculations rely primarily on the following expressions for the full conditional distributions:

$$\begin{aligned} \log p(\boldsymbol{\beta}, \mathbf{u}|\text{rest}) &= -\frac{1}{2} \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \left(\mathbf{C}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma})^{-1} \mathbf{C} + \begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\sigma_\ell^2 \mathbf{I}_{K_\ell})_{1 \leq \ell \leq d} \end{bmatrix} \right) \right. \\ &\quad \left. \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} - 2 \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \mathbf{C}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma})^{-1} \mathbf{y} \right) + \text{const}, \\ \log p(\sigma_\ell^2|\text{rest}) &= \left(\frac{1}{2} (K_\ell + 1) - 1 \right) \log(\sigma_\ell^2) - \left(\frac{1}{2} \|\mathbf{u}_\ell\|^2 + a_\ell^{-1} \right) / \sigma_\ell^2 + \text{const}, \\ \log p(a_\ell|\text{rest}) &= -2 \log(a_\ell) - (\sigma_\ell^2 + A^{-2}) / a_\ell + \text{const}, \\ \log p(\boldsymbol{\Sigma}|\text{rest}) &= -\frac{A_{\boldsymbol{\Sigma}} + m + n + 1}{2} \log|\boldsymbol{\Sigma}| \\ &\quad - \frac{1}{2} \text{tr} \left\{ \left[\sum_{i=1}^m (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{u}) (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{u})^\top + B_{\boldsymbol{\Sigma}} \right] \boldsymbol{\Sigma}^{-1} \right\} + \text{const}. \end{aligned}$$

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

We begin with

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}) &= E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u}|\text{rest}) \} + \text{const} \\ &= -\frac{1}{2} \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \left(\mathbf{C}^\top (\mathbf{I}_m \otimes M_{q(\boldsymbol{\Sigma}^{-1})}) \mathbf{C} \right. \right. \\ &\quad \left. \left. + \begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\mu_{q(1/\sigma_\ell^2)} \mathbf{I}_{K_\ell})_{1 \leq \ell \leq d} \end{bmatrix} \right) \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right. \\ &\quad \left. - 2 \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \mathbf{C}^\top (\mathbf{I}_m \otimes M_{q(\boldsymbol{\Sigma}^{-1})}) \mathbf{y} \right) + \text{const}. \end{aligned}$$

Therefore,

$q^*(\boldsymbol{\beta}, \mathbf{u})$ is the $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ density function,

where

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

and

$$\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} = \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \left(\mathbf{I}_m \otimes M_{q(\boldsymbol{\Sigma}^{-1})} \right) \mathbf{y}.$$

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

Next, we have

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

Therefore,

$$q^*(\sigma_\ell^2) \text{ is the Inverse-Gamma } \left(\frac{1}{2} (K_\ell + 1), B_{q(\sigma_\ell^2)} \right) \text{ density function,}$$

where

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

In addition, using Result 1.4.3, this gives

$$\mu_{q(1/\sigma_\ell^2)} = \frac{1}{2} (K_\ell + 1) / B_{q(\sigma_\ell^2)}.$$

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

Similarly,

$$\begin{aligned} \log q^*(a_\ell) &= E_q \{ \log p(a_\ell | \text{rest}) \} + \text{const} \\ &= -2 \log(a_\ell) - E_q(\sigma_\ell^2 + A^{-2}) / a_\ell + \text{const} \\ &= (-1 - 1) \log(a_\ell) - \left(\mu_{q(1/\sigma_\ell^2)} + A^{-2} \right) / a_\ell + \text{const}. \end{aligned}$$

Therefore

$$q^*(a_\ell) \text{ is the Inverse-Gamma } (1, B_{q(a_\ell)}) \text{ density function,}$$

where

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

Also, making use of Result 1.4.3, we get

$$\mu_{q(1/a_\ell)} = 1 / B_{q(a_\ell)}.$$

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

Finally,

$$\begin{aligned}
 \log q^*(\Sigma) &= E_q \{ \log p(\Sigma | \text{rest}) \} + \text{const} \\
 &= -\frac{1}{2} (A_\Sigma + m + n + 1) \log |\Sigma| \\
 &\quad - \frac{1}{2} \text{tr} \left[\sum_{i=1}^m E_{q(\beta, \mathbf{u})} \{ (\mathbf{y}_i - \mathbf{X}_i \beta - \mathbf{Z}_i \mathbf{u}) (\mathbf{y}_i - \mathbf{X}_i \beta - \mathbf{Z}_i \mathbf{u})^\top + B_\Sigma \} \Sigma^{-1} \right] \\
 &\quad + \text{const} \\
 &= -\frac{1}{2} (A_\Sigma + m + n + 1) \log |\Sigma| \\
 &\quad - \frac{1}{2} \text{tr} \left(\left[\sum_{i=1}^m \{ (\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\beta, \mathbf{u})}) (\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\beta, \mathbf{u})})^\top + \mathbf{C}_i \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \mathbf{C}_i^\top \} \right. \right. \\
 &\quad \left. \left. + B_\Sigma \right] \Sigma^{-1} \right) + \text{const}.
 \end{aligned}$$

Therefore,

$q^*(\Sigma)$ is the Inverse-Wishart $(A_\Sigma + m, B_{q(\Sigma)})$ density function,

where

$$B_{q(\Sigma)} = B_\Sigma + \sum_{i=1}^m \left[(\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\beta, \mathbf{u})}) (\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\beta, \mathbf{u})})^\top + \mathbf{C}_i \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \mathbf{C}_i^\top \right].$$

Lastly, making use of Result 1.4.7 gives

$$M_{q(\Sigma^{-1})} = (A_\Sigma + m) B_{q(\Sigma)}^{-1}.$$

2.B Derivation for the marginal log-likelihood lower bound

The expression for the lower bound on the marginal log-likelihood given at (2.17) is

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q) &= \frac{1}{2} \left(\sum_{\ell=1}^d K_\ell + p \right) - \frac{m}{2} \log(2\pi) - \log(\pi) + \sum_{\ell=1}^d \log \Gamma\left(\frac{1}{2}(K_\ell + 1)\right) \\
&\quad - \frac{p}{2} \log(\sigma_\beta^2) - \log(A) - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| \\
&\quad - \frac{1}{2} \sum_{\ell=1}^d (K_\ell + 1) \log(B_{q(\sigma_\ell^2)}) - \frac{1}{2} (A_\Sigma + m) \log |B_{q(\Sigma)}| + \sum_{\ell=1}^d \mu_{q(1/\sigma_\ell^2)} \mu_{q(1/a_\ell)} \\
&\quad - \sum_{\ell=1}^d \log(\mu_{q(1/\sigma_\ell^2)} + A^{-2}) + \frac{1}{2} A_\Sigma (\log |B_\ell|) - \log(C_{n,A}) + \log(C_{n,A+m}).
\end{aligned}$$

Derivation: The lower bound on the marginal log-likelihood is achieved through the following expression:

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

First we note that

$$\begin{aligned}
\log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}) &= -\frac{m}{2} \log(2\pi) - \frac{m}{2} \log |\boldsymbol{\Sigma}| \\
&\quad - \frac{1}{2} \text{tr} \{ (\mathbf{I}_m \otimes \boldsymbol{\Sigma})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u})^\top \}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
E_q \{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\Sigma}) \} &= -\frac{m}{2} \log(2\pi) - \frac{m}{2} E_q (\log |\boldsymbol{\Sigma}|) \\
&\quad - \frac{1}{2} \text{tr} \left\{ \left(\mathbf{I}_m \otimes M_{q(\boldsymbol{\Sigma}^{-1})} \right) \left[\mathbf{C} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \mathbf{C}^\top \right. \right. \\
&\quad \left. \left. + (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}) (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})})^\top \right] \right\}.
\end{aligned}$$

Next,

$$\begin{aligned}
&E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\sigma}^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}) \} \\
&= -\frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{2} \sum_{\ell=1}^d K_\ell E_q \{ \log(\sigma_\ell^2) \} - \frac{1}{2} \sum_{\ell=1}^d E_q(1/\sigma_\ell^2) E_q(\|\mathbf{u}_\ell\|^2) + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| \\
&\quad + \frac{1}{2} \left(\sum_{\ell=1}^d K_\ell + p \right)
\end{aligned}$$

2.B. DERIVATION FOR THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

$$\begin{aligned}
&= -\frac{1}{2}p \log(\sigma_\beta^2) - \frac{1}{2} \sum_{\ell=1}^d K_\ell E_q(\log(\sigma_\ell^2)) - \frac{1}{2\sigma_\beta^2} \{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \} \\
&\quad - \frac{1}{2} \sum_{\ell=1}^d \mu_{q(1/\sigma_\ell^2)} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}_\ell)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell)}) \} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| + \frac{1}{2} \left(\sum_{\ell=1}^d K_\ell + p \right).
\end{aligned}$$

Next,

$$\begin{aligned}
&\log p(\boldsymbol{\sigma}^2 | \mathbf{a}) - \log q^*(\boldsymbol{\sigma}^2) \\
&= \log \left[\prod_{\ell=1}^d \frac{\frac{1}{a_\ell}}{\Gamma(\frac{1}{2})} (\sigma_\ell^2)^{-\frac{1}{2}-1} \exp\left(-\frac{1}{\sigma_\ell^2}\right) \right] \\
&\quad - \log \left[\prod_{\ell=1}^d \frac{B_{q(\sigma_\ell^2)}^{\frac{1}{2}(K_\ell+1)}}{\Gamma(\frac{1}{2}(K_\ell+1))} (\sigma_\ell^2)^{-\frac{1}{2}(K_\ell+1)-1} \exp\left(\frac{-B_{q(\sigma_\ell^2)}}{\sigma_\ell^2}\right) \right] \\
&= -\frac{1}{2} \sum_{\ell=1}^d \log(a_\ell) - \log \Gamma(\frac{1}{2}) - \frac{3}{2} \sum_{\ell=1}^d \log(\sigma_\ell^2) - \sum_{\ell=1}^d a_\ell^{-1} (1/\sigma_\ell^2) + \sum_{\ell=1}^d B_{q(\sigma_\ell^2)}/\sigma_\ell^2 \\
&\quad + \sum_{\ell=1}^d \left(\frac{1}{2}(K_\ell+1) - 1 \right) \log(\sigma_\ell^2) + \sum_{\ell=1}^d \log \Gamma(\frac{1}{2}(K_\ell+1)) - \frac{1}{2} \sum_{\ell=1}^d (K_\ell+1) \log(B_{q(\sigma_\ell^2)}) \\
&= -\frac{1}{2} \sum_{\ell=1}^d \log(a_\ell) + \sum_{\ell=1}^d \log \Gamma(\frac{1}{2}(K_\ell+1)) - \frac{1}{2} \log(\pi) + \sum_{\ell=1}^d (B_{q(\sigma_\ell^2)} - a_\ell^{-1}) (1/\sigma_\ell^2) \\
&\quad + \frac{1}{2} \sum_{\ell=1}^d K_\ell \log(\sigma_\ell^2) - \frac{1}{2} \sum_{\ell=1}^d (K_\ell+1) \log(B_{q(\sigma_\ell^2)}).
\end{aligned}$$

This means that

$$\begin{aligned}
E_q \{ \log p(\boldsymbol{\sigma}^2 | \mathbf{a}) - \log q^*(\boldsymbol{\sigma}^2) \} &= -\frac{1}{2} \sum_{\ell=1}^d E_q \{ \log(a_\ell) \} + \sum_{\ell=1}^d \log \Gamma\left\{ \frac{1}{2}(K_\ell+1) \right\} \\
&\quad - \frac{1}{2} \log(\pi) + \sum_{\ell=1}^d (B_{q(\sigma_\ell^2)} - \mu_{q(1/a_\ell)}) \mu_{q(1/\sigma_\ell^2)} \\
&\quad + \frac{1}{2} \sum_{\ell=1}^d K_\ell E_q \{ \log(\sigma_\ell^2) \} - \frac{1}{2} \sum_{\ell=1}^d (K_\ell+1) \log(B_{q(\sigma_\ell^2)}).
\end{aligned}$$

In addition,

$$\begin{aligned}
&\log p(\mathbf{a}) - \log q^*(\mathbf{a}) \\
&= -\log(A) - \frac{1}{2} \log(\pi) + \frac{1}{2} \sum_{\ell=1}^d \log(a_\ell) - \sum_{\ell=1}^d \log(B_{q(a_\ell)}) + \sum_{\ell=1}^d (B_{q(a_\ell)} - A^{-2}) (1/a_\ell) \\
&= -\log(A) - \frac{1}{2} \log(\pi) + \frac{1}{2} \sum_{\ell=1}^d \log(a_\ell) - \sum_{\ell=1}^d \log(\mu_{q(1/\sigma_\ell^2)} + A^{-2}) + \sum_{\ell=1}^d \mu_{q(1/\sigma_\ell^2)} / a_\ell
\end{aligned}$$

2.B. DERIVATION FOR THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

and so

$$\begin{aligned} E_q \{ \log p(\mathbf{a}) - \log q^*(\mathbf{a}) \} &= -\log(A) - \frac{1}{2} \log(\pi) + \frac{1}{2} \sum_{\ell=1}^d E_q(\log(\mathbf{a}_\ell)) \\ &\quad - \sum_{\ell=1}^d \log\left(\mu_{q(1/\sigma_\ell^2)} + A^{-2}\right) + \sum_{\ell=1}^d \mu_{q(1/\sigma_\ell^2)} \mu_{q(1/a_\ell)}. \end{aligned}$$

Next,

$$\begin{aligned} \log p(\boldsymbol{\Sigma}) - \log q^*(\boldsymbol{\Sigma}) &= -\log(C_{n,A}) + \log(C_{n,A+m}) + \frac{m}{2} \log|\boldsymbol{\Sigma}| + \frac{1}{2} A \log|B| \\ &\quad - \frac{1}{2}(A+m) \log|B_{q(\boldsymbol{\Sigma})}| - \frac{1}{2} \text{tr}(B\boldsymbol{\Sigma}^{-1}) + \frac{1}{2} \text{tr}(B_{q(\boldsymbol{\Sigma})}\boldsymbol{\Sigma}^{-1}). \end{aligned}$$

Hence,

$$\begin{aligned} E_q \{ \log p(\boldsymbol{\Sigma}) - \log q^*(\boldsymbol{\Sigma}) \} &= -\log(C_{n,A}) + \log(C_{n,A+m}) + \frac{m}{2} E_q(\log|\boldsymbol{\Sigma}|) + \frac{1}{2} A \log|B| \\ &\quad - \frac{1}{2}(A+m) \log|B_{q(\boldsymbol{\Sigma})}| + \frac{1}{2} \text{tr}\left((B_{q(\boldsymbol{\Sigma})} - B) M_{q(\boldsymbol{\Sigma}^{-1})}\right). \end{aligned}$$

Next, note that

$$\begin{aligned} \sum_{\ell=1}^d \left(B_{q(\sigma_\ell^2)} - \mu_{q(1/a_\ell)} \right) \mu_{q(1/\sigma_\ell^2)} &= \sum_{\ell=1}^d \left[\frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}_\ell)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell)}) \right\} + \mu_{q(1/a_\ell)} \right. \\ &\quad \left. - \mu_{q(1/a_\ell)} \right] \mu_{q(1/\sigma_\ell^2)} \\ &= \frac{1}{2} \sum_{\ell=1}^d \mu_{q(1/\sigma_\ell^2)} \left(\|\boldsymbol{\mu}_{q(\mathbf{u}_\ell)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}_\ell)}) \right). \end{aligned}$$

We also note that

$$\begin{aligned} &\frac{1}{2} \text{tr} \left\{ (B_{q(\boldsymbol{\Sigma})} - B) M_{q(\boldsymbol{\Sigma}^{-1})} \right\} \\ &= \frac{1}{2} \text{tr} \left(\left[\sum_{i=1}^m \left\{ (\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) (\mathbf{y}_i - \mathbf{C}_i \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top + \mathbf{C}_i \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}_i^\top \right\} \right] M_{q(\boldsymbol{\Sigma}^{-1})} \right) \\ &= \frac{1}{2} \text{tr} \left[\left(\mathbf{I}_m \otimes M_{q(\boldsymbol{\Sigma}^{-1})} \right) \left\{ (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) (\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^\top + \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \right\} \right] \end{aligned}$$

The cancellations lead by these equalities leads to the lower bound expression given in (2.17).

2.C Derivation of the Fisher information matrix for the linear marginal longitudinal model

We derive the Fisher information matrix for the frequentist Gaussian marginal longitudinal regression model:

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}).$$

The log-likelihood of $(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = -\frac{nm}{2} \log(2\pi) - \frac{m}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

To begin, with the help of Result 1.4.1 (g), we take the first differential of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ with respect to $\boldsymbol{\beta}$:

$$\begin{aligned} d_{\boldsymbol{\beta}}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \frac{1}{2}(-\mathbf{X}d\boldsymbol{\beta})^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X}d\boldsymbol{\beta} \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X}d\boldsymbol{\beta} \end{aligned}$$

and then the second differential of $\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ with respect to $\boldsymbol{\beta}$ is

$$\begin{aligned} d_{\boldsymbol{\beta}}^2\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= (-\mathbf{X}d\boldsymbol{\beta})^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X}d\boldsymbol{\beta} \\ &= -(d\boldsymbol{\beta})^\top \mathbf{X}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X} (d\boldsymbol{\beta}), \end{aligned}$$

so, from Theorem 1.4.2, we see that

$$\mathbf{H}_{\boldsymbol{\beta}}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = -\mathbf{X}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X},$$

and so, the $\boldsymbol{\beta}$ block of $\mathbf{I}(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ is

$$-\mathbf{E}\{\mathbf{H}_{\boldsymbol{\beta}}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})\} = \mathbf{X}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\mathbf{X}.$$

In order to check that the maximum likelihood estimates of $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ are indeed uncorrelated, we next deal with

$$d_{\text{vec}(\boldsymbol{\Sigma})}\{d_{\boldsymbol{\beta}}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})\} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top d\{(\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})\}d\boldsymbol{\beta}.$$

However, we see that

$$\mathbf{E}[d_{\text{vec}(\boldsymbol{\Sigma})}\{d_{\boldsymbol{\beta}}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})\}] = 0.$$

2.C. DERIVATION OF THE FISHER INFORMATION MATRIX FOR THE LINEAR MARGINAL LONGITUDINAL MODEL

It then follows that $\mathbf{I}(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ is of block-diagonal form. We next move on to find the $\boldsymbol{\Sigma}$ block of the information matrix.

$$d_{\text{vec}(\boldsymbol{\Sigma})}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = -\frac{m}{2}d(\log|\boldsymbol{\Sigma}|) - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top d(\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Then, using Results 1.4.1 (c), (e) and (f), we get

$$d_{\text{vec}(\boldsymbol{\Sigma})}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = -\frac{m}{2}\text{tr}(\boldsymbol{\Sigma}^{-1}d\boldsymbol{\Sigma}) + \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \{ \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^{-1}d\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}) \} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

and so

$$\begin{aligned} d_{\text{vec}(\boldsymbol{\Sigma})}^2\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \frac{m}{2}\text{tr}\{ \boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma}) \} \\ &\quad - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \{ \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}) \} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &\quad - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \{ \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}) \} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \frac{m}{2}\text{tr}\{ \boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma}) \} \\ &\quad - (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \{ \mathbf{I}_m \otimes (\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}) \} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \end{aligned}$$

With the help of Results 1.4.18, 1.4.14 and 1.4.15, we see that

$$\begin{aligned} -\mathbb{E}\{d_{\text{vec}(\boldsymbol{\Sigma})}^2\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})\} &= -\frac{m}{2}\text{tr}\{ \boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma}) \} \\ &\quad + \text{tr}\{ \mathbf{I}_m \otimes \{ \boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1} \} \} \\ &= -\frac{m}{2}\text{tr}\{ \boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma}) \} \\ &\quad + m\text{tr}\{ \boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1} \} \\ &= \frac{m}{2}\text{tr}\{ \boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma}) \}. \end{aligned}$$

Then, using Result 1.4.16, we get

$$-\mathbb{E}\{d_{\text{vec}(\boldsymbol{\Sigma})}^2\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma})\} = \frac{m}{2}\text{vec}(d\boldsymbol{\Sigma})^\top (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(d\boldsymbol{\Sigma}).$$

Therefore,

$$\mathbf{H}_{\text{vec}(\boldsymbol{\Sigma})}\ell(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = -\frac{m}{2}(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}),$$

and the full information matrix is:

$$\mathbf{I}(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \begin{bmatrix} \mathbf{X}^\top (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{m}{2}(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \end{bmatrix}.$$

2.D Proof of Theorem 2.3.1

Theorem 3.4.7 (a) of Mardia *et al.* (1979) states that if

$$\mathbf{X} \sim \text{Wishart}(a, \mathbf{B}), \quad a > n, \quad \mathbf{B} \text{ positive definite}, \quad (2.21)$$

then the ratio $\mathbf{c}^\top \mathbf{B} \mathbf{c} / \mathbf{c}^\top \mathbf{X}^{-1} \mathbf{c}$ has a χ_{a-n+1}^2 distribution for any fixed $n \times 1$ vector \mathbf{c} . If we let $\mathbf{c} = [0 \ 1 \ 0]^\top$, such that $\mathbf{B}_{22} = \mathbf{c}^\top \mathbf{B} \mathbf{c}$ and $(\mathbf{X}^{-1})_{22} = \mathbf{c}^\top \mathbf{X}^{-1} \mathbf{c}$. Then,

$$\frac{\mathbf{B}_{22}}{(\mathbf{X}^{-1})_{22}} \sim \chi_{a-3+1}^2$$

and making use of Result 1.4.5, this is equivalent to

$$\frac{\mathbf{B}_{22}}{(\mathbf{X}^{-1})_{22}} \sim \text{Gamma}\left(\frac{a-3+1}{2}, \frac{1}{2}\right).$$

Using Result 1.4.4 we see that:

$$\frac{(\mathbf{X}^{-1})_{22}}{\mathbf{B}_{22}} \sim \text{Inverse-Gamma}\left(\frac{a-3+1}{2}, \frac{1}{2}\right).$$

Therefore using Results 1.4.2 and 1.4.8, the $[2, 2]$ entry of an Inverse-Wishart random matrix \mathbf{X}^{-1} has distribution

$$(\mathbf{X}^{-1})_{22} \sim \text{Inverse-Gamma}\left(\frac{a-3+1}{2}, \frac{1}{2} \mathbf{B}_{22}\right).$$

The distribution of the $[1, 1]$ and $[3, 3]$ entries of \mathbf{X}^{-1} have similar forms. We have shown this for a 3×3 random matrix, however this applies for general $n \times n$ random matrices also. Thus, for a general $n \times n$ positive definite matrix \mathbf{X} that is distributed

$$\mathbf{X} \sim \text{Inverse-Wishart}(a, \mathbf{B}),$$

its diagonal entries have the following distribution:

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

2.E Proof of Theorem 2.3.2

The characteristic function of a Wishart random matrix \mathbf{X} with distribution

$$\mathbf{X} \sim \text{Inverse-Wishart}(a, \mathbf{B})$$

is

$$E[\exp\{i \operatorname{tr}(\mathbf{X}\mathbf{T})\}] = |\mathbf{I} - 2i\mathbf{T}\mathbf{B}^{-1}|^{-a/2}.$$

If we let \mathbf{X} be of dimension 3×3 and focus on finding the characteristic function of the $[1, 2]$ and $[2, 1]$ off-diagonal entries of \mathbf{X} , we set

$$\mathbf{T} = \begin{bmatrix} 0 & t & 0 \\ t & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

so that $\operatorname{tr}(\mathbf{X}\mathbf{T}) = 2t\mathbf{X}_{12}$. Therefore the characteristic function becomes

$$\begin{aligned} E\{\exp(2it\mathbf{X}_{12})\} &= |\mathbf{I} - 2i\mathbf{T}\mathbf{B}^{-1}|^{-a/2} \\ &= \exp\left\{-\frac{a}{2} \log |\mathbf{I} - 2i\mathbf{T}\mathbf{B}^{-1}|\right\}, \end{aligned}$$

and since the determinant of a matrix is equal to the product of the eigenvalues of that matrix, we get

$$E\{\exp(2it\mathbf{X}_{12})\} = \exp\left\{-\frac{a}{2} \sum_{\ell=1}^n \log(\lambda_{\ell})\right\},$$

where the $\lambda_{\ell}, 1 \leq \ell \leq n$, are the eigenvalues of the matrix $|\mathbf{I} - 2i\mathbf{T}\mathbf{B}^{-1}|$. This can also be shown for Wishart random matrices of dimension $n \times n$.

Chapter 3

Variational inference for heteroscedastic semiparametric regression

3.1 Introduction

In many situations considering real data, the assumption of constant conditional variance on the regressors is erroneous, albeit used routinely in parametric and nonparametric regression. Ignorance of significant heteroscedasticity will lead to misguided inferential statistics such as prediction intervals. A remedy for this is to employ heteroscedastic nonparametric regression which overcomes these complications. For a set of data, (x_i, y_i) , $1 \leq i \leq n$, we replace the homoscedastic nonparametric regression model

$$E(y_i) = f(x_i), \quad \text{Var}(y_i) = \sigma^2$$

with

$$E(y_i) = f(x_i), \quad \text{Var}(y_i) = g(x_i), \tag{3.1}$$

where the variance function g is estimated simultaneously with the mean function f . Several approaches already exist for fitting (3.1), some of which are encompassed in Ruppert *et al.* (2003) and Rigby & Stasinopoulos (2005). However, graphical model representations of mixed model-based penalized splines (Wand, 2009) have shown to be an effective

The content of this chapter is published as: Menictas, M and Wand. M.P. (2015). Variational inference for heteroscedastic semiparametric regression. *Australian and New Zealand Journal of Statistics*, **57**, Number 1, 119–138. This research was also presented at the *Australian Statistical Conference in conjunction with the Institute of Mathematical Statistics Annual Meeting*, Sydney, 2014.

approach, especially when considering the extendability of models arising from (3.1).

In this chapter we investigate a relatively new modification of MFVB known as *non-conjugate variational message passing* (Knowles & Minka, 2011), which aids in the adaptation of heteroscedasticity to the appropriate model whilst also involving only closed form algebraic expressions. Even though we focus on univariate nonparametric regression, the modularity of our approach allows straightforward extensions to more complex models. This is known as the *locality* property of MFVB and is described in section 1.8.1 and in greater detail in section 3 of Wand *et al.* (2011). For example, if a complex graphical model includes a component where one variable is modelled as a heteroscedastic nonparametric regression function of another variable, the variational inference for the parameters in that section of the graph can be developed using the methodology in this chapter.

Lázaro-Gredilla & Titsias (2011) and Nott *et al.* (2012) have also developed variational inference for heteroscedastic regression models, although the latter contribution was restricted to linear mean and log-variance functions. Simultaneous nonparametric mean and variance function estimation aided by an elegant Gaussian process approach was achieved by Lázaro-Gredilla & Titsias (2011). Their strategy involved full-rank nonparametric function estimators with Gauss-Hermite quadrature for variance function estimation. Bugbee *et al.* (2015) use MFVB with an embedded Laplace approximation to account for semi-parametric regression with heteroscedasticity. Our approach is different by using low-rank penalized splines, as well as a non-conjugate MFVB algorithm that involves closed form updates, making it more susceptible to increasingly larger data sets.

In Section 3.2, we provide details on the Bayesian penalized spline model for simultaneous mean and variance function estimation based on univariate data. Section 3.3 provides a description of the variational inference methodology used to formulate Algorithm 2. In Sections 3.4 and 3.5 we present numerical analyses and results which confirm the speed achieved by MFVB.

3.2 Model description

The generic form of the Gaussian heteroscedastic nonparametric regression model is of the form

$$y_i \stackrel{\text{ind.}}{\sim} \text{N}(f(x_i), g(x_i)), \quad 1 \leq i \leq n, \quad (3.2)$$

3.2. MODEL DESCRIPTION

where (x_i, y_i) represents the i th predictor/response pair of a regression data-set and the functions f and g are real valued smooth functions. The quantity g is referred to as the variance function, and \sqrt{g} is the standard deviation function. Fitting of (3.2) using mixed model-based penalized splines (e.g. Ruppert *et al.*, 2003) requires the following structural forms of f and g :

$$\begin{aligned} f(x) &= \beta_0 + \beta_1 x + \sum_{k=1}^{K_u} u_k z_k^u(x), & u_k &\stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2) \\ \text{and } g(x) &= \exp\left(\gamma_0 + \gamma_1 x + \sum_{k=1}^{K_v} v_k z_k^v(x)\right), & v_k &\stackrel{\text{ind.}}{\sim} N(0, \sigma_v^2). \end{aligned} \quad (3.3)$$

The z_k^u , $1 \leq k \leq K_u$ and z_k^v , $1 \leq k \leq K_v$ are spline bases of sizes K_u and K_v respectively. Here we use suitably transformed cubic O'Sullivan splines as our default for the z_k^u and z_k^v , as explained in Section 4 of Wand & Ormerod (2008). Further, we allow for K_u and K_v to comprise different sizes. The prior specifications on the model parameters are

$$\beta_0, \beta_1 \stackrel{\text{ind.}}{\sim} N(0, \sigma_\beta^2), \quad \gamma_0, \gamma_1 \stackrel{\text{ind.}}{\sim} N(0, \sigma_\gamma^2), \quad \sigma_u \sim \text{Half-Cauchy}(A_u) \text{ and } \sigma_v \sim \text{Half-Cauchy}(A_v),$$

where $\sigma_\beta, \sigma_\gamma, A_u, A_v > 0$ are hyperparameters to be specified by the user. Under the assumption that the data has been transformed to have zero mean and unit variance, the default setting we employ for the hyperparameters throughout this chapter are:

$$\sigma_\beta = \sigma_\gamma = A_u = A_v = 10^5. \quad (3.4)$$

The complete Bayesian hierarchical model corresponding to (3.2) is

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{u}, \mathbf{v} &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_u \mathbf{u}, \text{diag}\{\exp(\mathbf{X}\boldsymbol{\gamma} + \mathbf{Z}_v \mathbf{v})\}), \\ \mathbf{u} | \sigma_u^2 &\sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}_{K_u}), \quad \mathbf{v} | \sigma_v^2 \sim N(\mathbf{0}, \sigma_v^2 \mathbf{I}_{K_v}), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_2), \quad \boldsymbol{\gamma} \sim N(\mathbf{0}, \sigma_\gamma^2 \mathbf{I}_2), \\ \sigma_u^2 | a_u &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_u), \quad a_u \sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_u^2), \\ \sigma_v^2 | a_v &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a_v), \quad a_v \sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A_v^2), \end{aligned} \quad (3.5)$$

where $\boldsymbol{\beta}$ is the 2×1 vector of fixed effects containing (β_0, β_1) and $\boldsymbol{\gamma}$ defined similarly. Further, \mathbf{u} is the $K_u \times 1$ vector containing (u_1, \dots, u_{K_u}) , \mathbf{v} is a $K_v \times 1$ vector defined similarly, and σ_u^2 and σ_v^2 are variance components corresponding to \mathbf{u} and \mathbf{v} respectively. The design matrices, \mathbf{X} , \mathbf{Z}_u and \mathbf{Z}_v , are defined to be

$$\mathbf{X} \equiv [1 \ x_i]_{1 \leq i \leq n}, \quad \mathbf{Z}_u \equiv \left[z_k^u(x_i) \right]_{\substack{1 \leq k \leq K_u \\ 1 \leq i \leq n}} \quad \text{and} \quad \mathbf{Z}_v \equiv \left[z_k^v(x_i) \right]_{\substack{1 \leq k \leq K_v \\ 1 \leq i \leq n}}.$$

3.3. VARIATIONAL INFERENCE ALGORITHM

It is beneficial to merge the mean function coefficients and variance function coefficients into single vectors:

$$\boldsymbol{\nu} \equiv \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \quad \text{and} \quad \boldsymbol{\omega} \equiv \begin{bmatrix} \boldsymbol{\gamma} \\ \mathbf{v} \end{bmatrix}. \quad (3.6)$$

Combining the corresponding design matrices $\mathbf{C}_\nu \equiv [\mathbf{X} \ \mathbf{Z}_u]$ and $\mathbf{C}_\omega \equiv [\mathbf{X} \ \mathbf{Z}_v]$ then allows us to use the following notation

$$\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_u\mathbf{u} = \mathbf{C}_\nu\boldsymbol{\nu} \quad \text{and} \quad \mathbf{X}\boldsymbol{\gamma} + \mathbf{Z}_v\mathbf{v} = \mathbf{C}_\omega\boldsymbol{\omega}.$$

The directed acyclic graph corresponding to the model shown in (3.5) is shown in Figure 3.1.

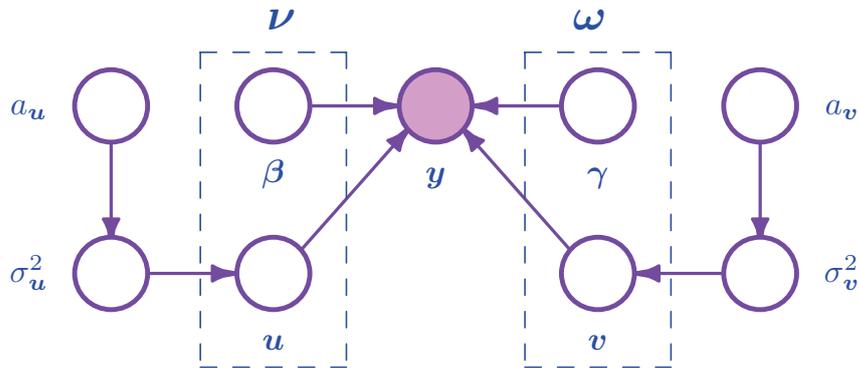


Figure 3.1: Directed acyclic graph for the model conveyed in (3.5). The shaded node corresponds to the observed data vector. Random effects and auxiliary variables are referred to as hidden nodes. The dashed boxes indicate that $\boldsymbol{\beta}$ and \mathbf{u} are combined into a vector denoted by $\boldsymbol{\nu}$ and $\boldsymbol{\omega}$ is the combination of $\boldsymbol{\gamma}$ and \mathbf{v} .

3.3 Variational inference algorithm

Having discussed the elementary design of our model, we next provide the methodology required to achieve variational bayesian inference for (3.5). Ordinary implementation of MFVB methodology involves restricting the full posterior density function to have the approximate product form

$$p(\boldsymbol{\nu}, \boldsymbol{\omega}, \mathbf{a}, \boldsymbol{\sigma}^2 | \mathbf{y}) \approx q(\boldsymbol{\nu})q(\boldsymbol{\omega})q(\sigma_u^2, \sigma_v^2)q(a_u, a_v), \quad (3.7)$$

where $\mathbf{a} = (a_u, a_v)$ and $\boldsymbol{\sigma}^2 = (\sigma_u^2, \sigma_v^2)$. Additional factorisations can easily be detected using moralisation. After moralising (see Definition 1.8.9) it becomes clear that all paths be-

3.3. VARIATIONAL INFERENCE ALGORITHM

tween σ_u^2 and σ_v^2 must visit at least one of the components in the set $\{\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}\}$. Using Theorem 1.8.1, σ_u^2 is separated from σ_v^2 given $\{\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}\}$. Similarly, $a_u \perp a_v | \{\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2\}$. Therefore, (3.7) reduces to

$$q(\boldsymbol{\nu})q(\boldsymbol{\omega})q(\sigma_u^2, \sigma_v^2)q(a_u, a_v) = q(\boldsymbol{\nu})q(\boldsymbol{\omega})q(\sigma_u^2)q(\sigma_v^2)q(a_u)q(a_v). \quad (3.8)$$

Under this product restriction, the parameters $\boldsymbol{\nu}, \sigma_u^2, \sigma_v^2, a_u$ and a_v all have closed form expressions for their full conditional distributions

$$\begin{aligned} p(\boldsymbol{\nu}|\text{rest}) &\sim N\left(\left\{\mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}}) \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & 0 \\ 0 & \sigma_u^{-2} \mathbf{I}_{K_u} \end{bmatrix}\right\} \mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}}) \mathbf{y}, \right. \\ &\quad \left. \left\{\mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}}) \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & 0 \\ 0 & \sigma_u^{-2} \mathbf{I}_{K_u} \end{bmatrix}\right\}^{-1}\right) \\ p(\sigma_u^2|\text{rest}) &\sim \text{Inverse-Gamma}\left(\frac{1}{2}(K_u + 1), \frac{1}{2}\|\mathbf{u}\|^2 + a_u^{-1}\right) \\ p(\sigma_v^2|\text{rest}) &\sim \text{Inverse-Gamma}\left(\frac{1}{2}(K_v + 1), \frac{1}{2}\|\mathbf{v}\|^2 + a_v^{-1}\right) \\ p(a_u) &\sim \text{Inverse-Gamma}\left(1, \sigma_u^{-2} + A_u^{-2}\right) \\ p(a_v) &\sim \text{Inverse-Gamma}\left(1, \sigma_v^{-2} + A_v^{-2}\right). \end{aligned}$$

However, the optimal posterior density function of $\boldsymbol{\omega}$ is

$$q^*(\boldsymbol{\omega}) \propto \exp\left[E_{q(-\boldsymbol{\omega})}\{\log p(\boldsymbol{\omega}|\text{rest})\}\right], \quad (3.9)$$

which involves intractable multivariate integrals that are not available in closed form expressions. A remedy in this instance is to work with the product restriction

$$p(\boldsymbol{\nu}, \boldsymbol{\omega}, \mathbf{a}, \boldsymbol{\sigma}^2 | \mathbf{y}) \approx q(\boldsymbol{\nu})q(\boldsymbol{\omega}; \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})q(\mathbf{a})q(\boldsymbol{\sigma}^2), \quad (3.10)$$

where

$$q(\boldsymbol{\omega}; \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) \text{ density function.} \quad (3.11)$$

Thus, the marginal log-likelihood is

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) &\equiv E_q\left[\log p(\boldsymbol{\nu}, \boldsymbol{\omega}, \mathbf{a}, \boldsymbol{\sigma}^2, \mathbf{y}) \right. \\ &\quad \left. - \log\{q(\boldsymbol{\nu})q(\boldsymbol{\omega}; \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})q(\mathbf{a})q(\boldsymbol{\sigma}^2)\}\right] \end{aligned} \quad (3.12)$$

and choosing $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ to maximise (3.12) corresponds to minimization of the Kullback-Leibler divergence. Knowles & Minka (2011) termed this approach *non-conjugate variational message passing*, since it offers a way of getting around the non-conjugacies of

3.3. VARIATIONAL INFERENCE ALGORITHM

ordinary MFVB. This approach is explained in more detail in Section 1.7.

For fixed $(\boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$, application of ordinary MFVB, with $q^*(\boldsymbol{\omega})$ omitted, gives

$$\begin{aligned}
q^*(\boldsymbol{\nu}) & \text{ is a Normal } (\boldsymbol{\mu}_{q(\boldsymbol{\nu})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}) \text{ density function,} \\
q^*(\sigma_u^2) & \text{ is an Inverse-Gamma } (\tfrac{1}{2}(K_u + 1), B_{q(\sigma_u^2)}) \text{ density function,} \\
q^*(\sigma_v^2) & \text{ is an Inverse-Gamma } (\tfrac{1}{2}(K_v + 1), B_{q(\sigma_v^2)}) \text{ density function,} \\
q^*(a_u) & \text{ is an Inverse-Gamma } (1, B_{q(a_u)}) \text{ density function} \\
\text{and } q^*(a_v) & \text{ is an Inverse-Gamma } (1, B_{q(a_v)}) \text{ density function,}
\end{aligned} \tag{3.13}$$

for parameters $\boldsymbol{\mu}_{q(\boldsymbol{\nu})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}$, the mean and covariance matrix of $q^*(\boldsymbol{\nu})$, $B_{q(\sigma_u^2)}$, the rate parameter of $q^*(\sigma_u^2)$, $B_{q(\sigma_v^2)}$, the rate parameter of $q^*(\sigma_v^2)$, $B_{q(a_u)}$, the rate parameter of $q^*(a_u)$ and $B_{q(a_v)}$, the rate parameter of $q^*(a_v)$.

With assistance from Knowles & Minka (2011) who propose a fixed-point iteration scheme for maximizing (3.12) and Wand (2014) who provides simplified fixed point updates for the Multivariate Normal q -density parameters, such as that shown in (3.11), we can now present Algorithm 2.

Convergence in Algorithm 2 is assessed using the variational lower bound on the approximate marginal log-likelihood, denoted by $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$ and has the following expression:

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) & = \tfrac{1}{2}(K_u + K_v + 4) - \tfrac{n}{2} \log(2\pi) + \log \Gamma(\tfrac{1}{2}(K_u + 1)) - 2 \log(\pi) \\
& + \log \Gamma(\tfrac{1}{2}(K_v + 1)) - \log(A_u) - \log(A_v) - \tfrac{1}{2} \mathbf{1}^T (\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})}) \\
& - \tfrac{1}{2} \mathbf{1}^T \{ \boldsymbol{\mu}_{q(\sigma_v^2)} \odot \exp(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \} - \log(\sigma_\beta^2) - \log(\sigma_\gamma^2) \\
& + \tfrac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}| + \tfrac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}| - \tfrac{1}{2\sigma_\beta^2} (\| \boldsymbol{\mu}_{q(\boldsymbol{\beta})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})})) \\
& - \tfrac{1}{2\sigma_\gamma^2} (\| \boldsymbol{\mu}_{q(\boldsymbol{\gamma})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\gamma})})) - \tfrac{1}{2}(K_u + 1) \log(B_{q(\sigma_u^2)}) \\
& - \tfrac{1}{2}(K_v + 1) \log(B_{q(\sigma_v^2)}) - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \\
& - \log(\mu_{q(1/\sigma_v^2)} + A_v^{-2}) + \mu_{q(1/\sigma_u^2)} \mu_{q(1/a_u)} + \mu_{q(1/\sigma_v^2)} \mu_{q(1/a_v)}.
\end{aligned}$$

The optimal parameters shown in Algorithm 2 are all interdependent and are obtained by cycling through this iterative coordinate ascent procedure. The algorithm also provides fixed point iterative updates for $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$. Details on the derivation of these parameters, as well as the parameters in (3.13) are given in Appendix 3.A. In addition, $\boldsymbol{\mu}_{q(\mathbf{u})}$ is defined to be the sub-vector of $\boldsymbol{\mu}_{q(\boldsymbol{\nu})}$ corresponding to \mathbf{u} , and $\boldsymbol{\Sigma}_{q(\mathbf{u})}$ is the sub-matrix of $\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}$ corresponding to \mathbf{u} . The parameters $\boldsymbol{\mu}_{q(\mathbf{v})}$ and $\boldsymbol{\Sigma}_{q(\mathbf{v})}$ are defined similarly. It is important to note that unlike ordinary MFVB, there is no guarantee that there will be an increase at each iteration for $\underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$. Therefore, the cycle should end

Set up initial values:

$\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ a $(K_v + 2) \times 1$ vector, $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ a $(K_v + 2) \times (K_v + 2)$ positive definite matrix, $\boldsymbol{\mu}_{q(1/\sigma_u^2)}, \boldsymbol{\mu}_{q(1/\sigma_v^2)} > 0$, and $\boldsymbol{\mu}_{q(r_\nu^2)}$ an $n \times 1$ vector.

Cycle through:

$$\boldsymbol{\psi}_{q(\boldsymbol{\omega})} \leftarrow \exp \left\{ -\mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2} \text{diagonal} \left(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}^\top \right) \right\}$$

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})} \leftarrow \left(\mathbf{C}^\top \text{diag} \{ \boldsymbol{\psi}_{q(\boldsymbol{\omega})} \} \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_u^2)} \mathbf{I}_K \end{bmatrix} \right)^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\nu})} \leftarrow \boldsymbol{\Sigma}_{q(\boldsymbol{\nu})} \mathbf{C}^\top \text{diag} \{ \boldsymbol{\psi}_{q(\boldsymbol{\omega})} \} \mathbf{y}$$

$$\boldsymbol{\mu}_{q(r_\nu^2)} \leftarrow \text{diagonal} \left\{ (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\mu}_{q(\boldsymbol{\nu})}) (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\mu}_{q(\boldsymbol{\nu})})^\top + \mathbf{C}_\nu \boldsymbol{\Sigma}_{q(\boldsymbol{\nu})} \mathbf{C}_\nu^\top \right\}$$

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \leftarrow \left(\mathbf{C}^\top \text{diag} \{ \boldsymbol{\mu}_{q(r_\nu^2)} \odot \boldsymbol{\psi}_{q(\boldsymbol{\omega})} \} \mathbf{C} + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I}_K \end{bmatrix} \right)^{-1}$$

$$\boldsymbol{\mu}_{q(\boldsymbol{\omega})} \leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \left\{ \mathbf{C}^\top (\boldsymbol{\mu}_{q(r_\nu^2)} \odot \boldsymbol{\psi}_{q(\boldsymbol{\omega})} - \mathbf{1}) - \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I}_K \end{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\omega})} \right\}$$

$$\mu_{q(1/a_u)} \leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^2) \quad ; \quad \mu_{q(1/a_v)} \leftarrow 1 / (\mu_{q(1/\sigma_v^2)} + A_v^2)$$

$$\mu_{q(1/\sigma_u^2)} \leftarrow (K_u + 1) / \{ 2\mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \}$$

$$\mu_{q(1/\sigma_v^2)} \leftarrow (K_v + 1) / \{ 2\mu_{q(1/a_v)} + \|\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})}) \}$$

until the absolute relative change in $\log p(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$ is negligible.

Algorithm 2: *MFVB algorithm for the determination of the optimal parameters in $q^*(\boldsymbol{\omega})$, $q^*(\boldsymbol{\nu})$, $q^*(\sigma_u^2)$, $q^*(\sigma_v^2)$, $q^*(a_u)$ and $q^*(a_v)$.*

3.4. ASSESSMENT OF PERFORMANCE

Setting	$f(x)$	$\log g(x)$
A	$\sin(3\pi x^2)$	$0.1 + \cos(4\pi x)$
B	$-1.02x + 0.018x^2 + 5\phi(x; 0.38, 0.08) + \frac{8}{3}\phi(x; 0.75, 0.03)$	$-0.5 - \Phi(x; 0.2, 0.1) + 0.3x^2$
C	$\frac{35}{8}\phi(x; 0.01, 0.08) + \frac{190}{23}\phi(x; 0.45, 0.23) + 1.8\left\{1 - \frac{100}{14}\phi(x; 0.7, 0.14)\right\}$	$\frac{3}{2}\phi(x; 0, 0.2) + 4\phi(x; 1, 0.1)$
D	$\sin(3\pi x^2) - 1.02x + 0.018x^2 + 5\phi(x; 0.38, 0.08)$	$\cos(4\pi x) - 0.4 + 0.3x^2 - \Phi(x; 0.2, 10)$

Table 3.1: *Details of simulation study settings.*

once the absolute relative change in its logarithm falls below a negligible amount.

Algorithm 2 has been applied to simulated data and accuracy against MCMC has been assessed. The results are summarized in Section 3.4.

3.4 Assessment of performance

In order to assess the performance of Algorithm 2, we carried out an extensive simulation study to address the accuracy and computing time of model (3.2). Data were simulated according to model (3.2) with $n = 500$ and the x_i s uniform on $(0, 1)$. The four mean and variance function pairs that were used are listed in Table 3.1, where $\phi(\cdot; \mu, \sigma)$ and $\Phi(\cdot; \mu, \sigma)$ denote the density and distribution functions of the Normal distribution with mean μ and standard deviation σ , respectively. 100 data-sets were generated for each of the function pairs shown in Table 3.1. Each model corresponding to a new replication was fitted using MFVB corresponding to Algorithm 2 and MCMC based on a burnin of 5000, kept sample of 5000 and thinning factor of 5. The MFVB iterations were terminated when the relative change in $\log \underline{p}(\mathbf{y}; q)$ fell below 10^{-7} .

In the sections that follow, we provide details on the accuracy, coverage and timing of MFVB against the MCMC benchmark. Let us next discuss each of these assessments in more detail.

3.4.1 Assessment of accuracy

Fast approximate inference for the model parameters is achieved by Algorithm 2, however there is no guarantee that an acceptable level of accuracy will occur. Figure 3.2 provides

3.4. ASSESSMENT OF PERFORMANCE

an accuracy assessment of Algorithm 2, where the accuracy scores are summarized by boxplots for each parameter of interest. The accuracy score is defined in Section 1.9.

The parameters of interest for which accuracy is monitored are the $f(H_k)$ and $g(H_k)$, $1 \leq k \leq 5$, where the H_k are the sample hexiles of the x_i s. It is evident from Figure 3.2 that most of the accuracies for the $f(H_k)$ lie around 90%, while accuracies for the $g(H_k)$ lie around 80%. These favourable accuracy results are in keeping with the heuristics given in Section 2.3.1 of Chapter 2.

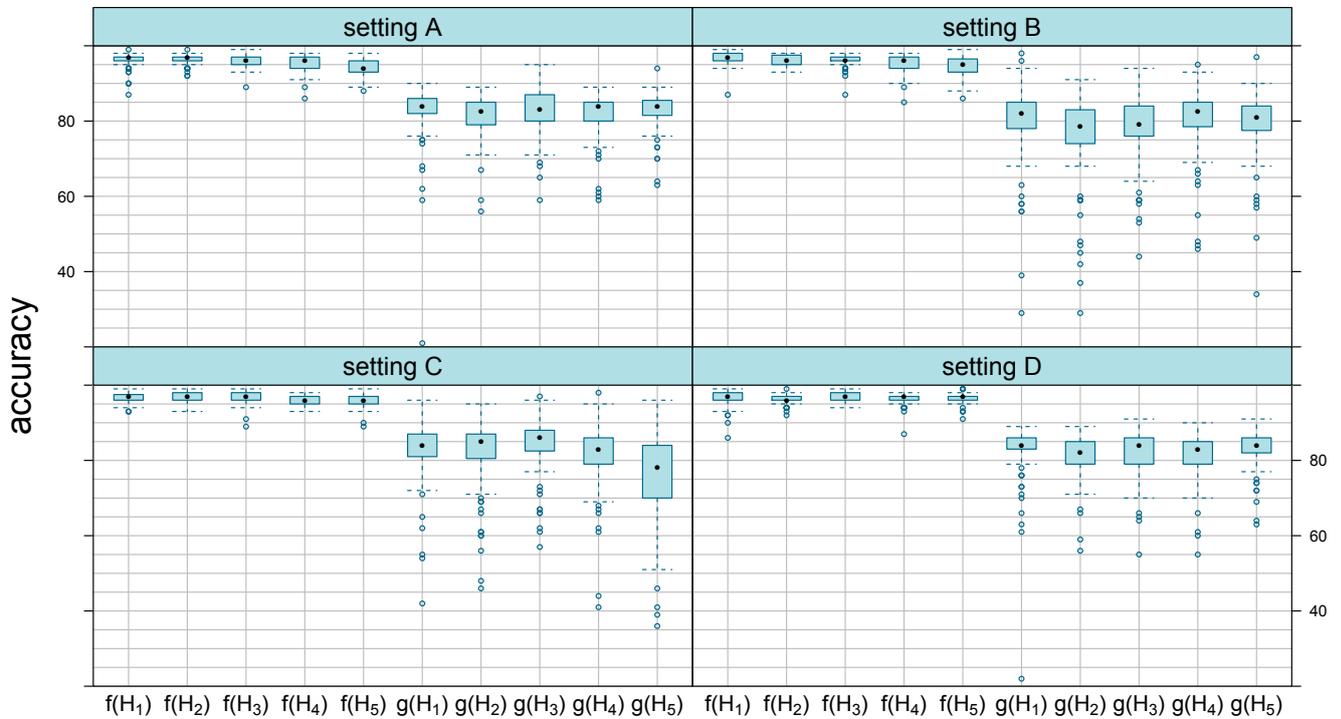


Figure 3.2: Summary of simulation study where the accuracy values are summarized using boxplots.

The fitted mean and standard deviation functions for the first replication in each of the four simulation settings is shown in Figure 3.3, with comparisons of MCMC and MFVB. As illustrated, the MCMC and MFVB fits show very good agreement for the mean functions and relatively good agreement for the standard deviation functions.

It is also of interest to see how well the joint distribution of certain parameters are doing in terms of approximation accuracy. Figure 3.4 shows a visual assessment of the approximation accuracy for the joint posterior density functions of the mean and variance hexile pairs according to simulation setting A.

3.4. ASSESSMENT OF PERFORMANCE

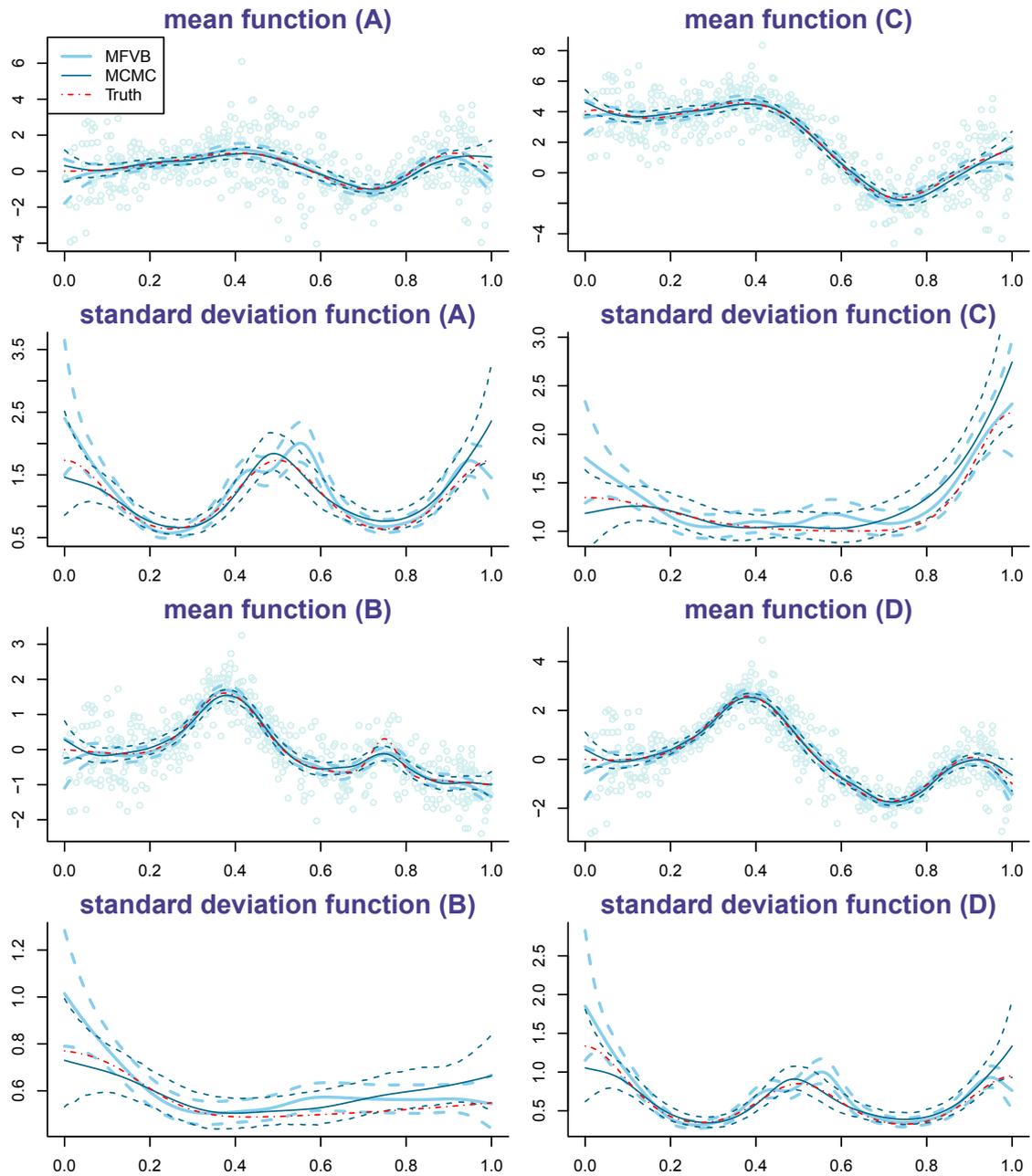


Figure 3.3: Comparison of MCMC and MFVB fitted functions for the first data-set in each of the four simulation settings.

In many cases, as shown by Figure 3.2, the accuracies can provide evidence of poor MFVB performance. Figure 3.5 illustrates some of the poorer approximations for this simulation study.

3.4. ASSESSMENT OF PERFORMANCE

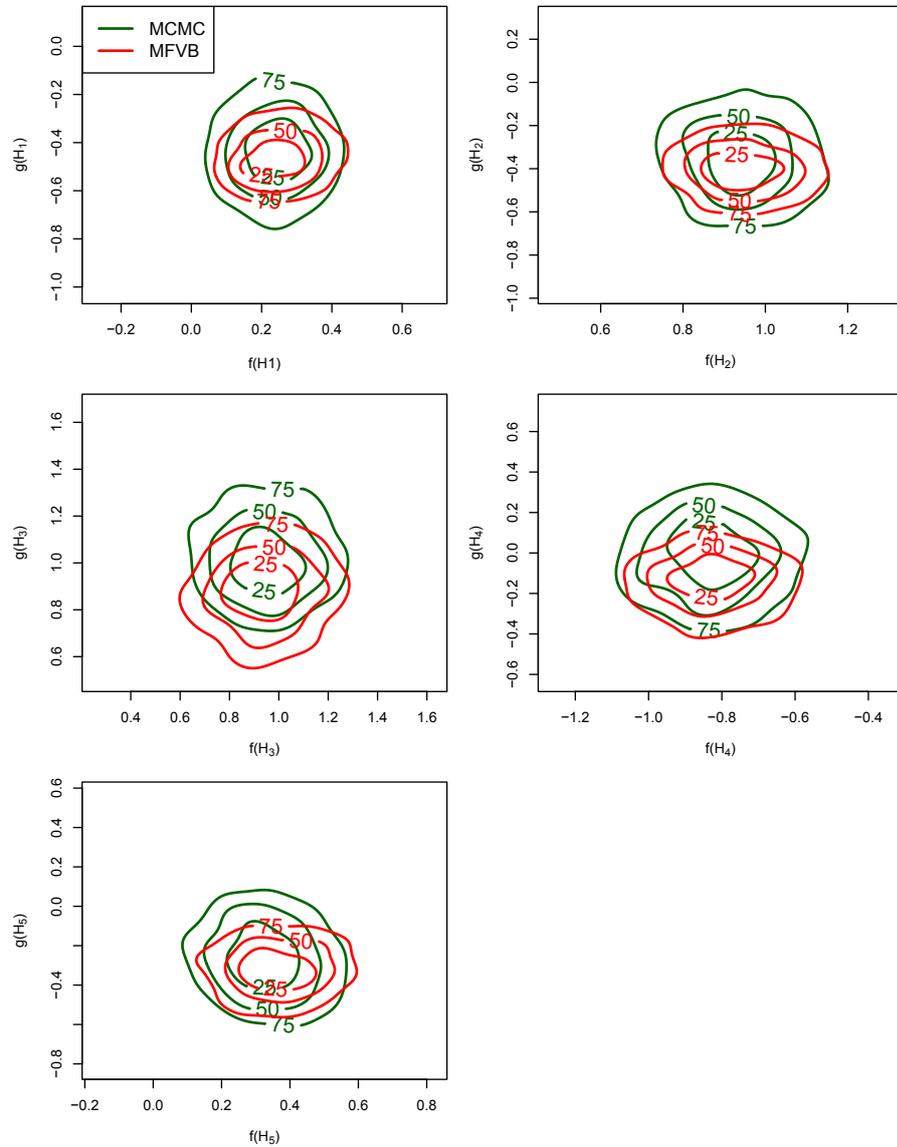


Figure 3.4: *Example of approximation accuracy for the joint posterior distributions of the mean and variance hexiles given for the first dataset in simulation setting A.*

3.4.2 Assessment of coverage

Another assessment of the performance of MFVB, involves the comparison between the true coverage and the coverage gained by the MFVB approximate credible intervals. The percentages of the true parameter coverage based on the approximate 95% credible intervals attained from the MFVB posterior densities are given in Table 3.2. As can be seen, the coverage overall is good and does not fall below 86%. Here we provide further confirmation that the performance of MFVB is excellent for the mean function and relatively

3.4. ASSESSMENT OF PERFORMANCE

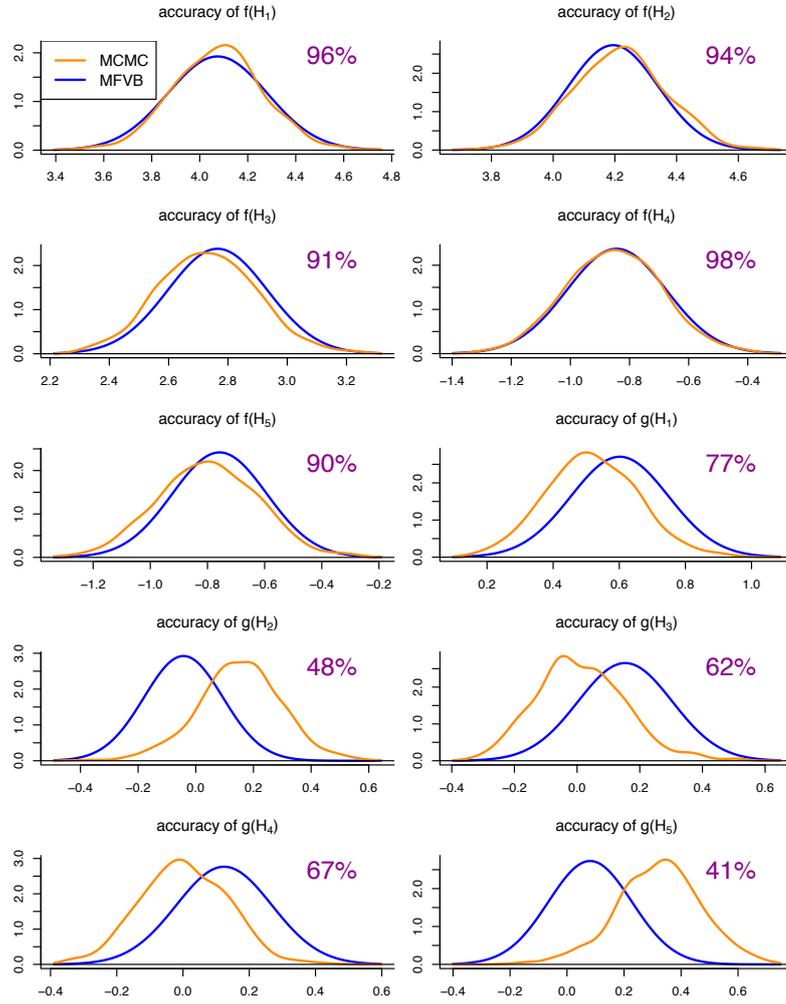


Figure 3.5: *Example of a poor approximation (for the variance function) for a dataset corresponding to simulation setting C.*

good for the variance function.

	$f(H_1)$	$f(H_2)$	$f(H_3)$	$f(H_4)$	$f(H_5)$	$g(H_1)$	$g(H_2)$	$g(H_3)$	$g(H_4)$	$g(H_5)$
sett. A	98	98	94	98	97	89	87	83	87	82
sett. B	98	98	95	96	95	83	83	90	89	83
sett. C	99	98	96	99	99	90	91	92	90	76
sett. D	98	98	95	98	98	89	89	83	86	82

Table 3.2: *Percentage coverage of the true parameter values by approximate 95% credible intervals based on variational Bayes approximate posterior density functions. The percentages are based on 100 replications.*

3.4.3 Assessment of speed

In order to demonstrate the time savings to be gained by using MFVB when compared to MCMC, we monitored the time taken to fit each model for each approach. The study was performed on a desktop computer (Intel Core i5-2400 3.10 GHz processor, 8 GBytes of random access memory). The resulting times are summarized in Table 3.3.

As we discussed previously, convergence for both approaches was assessed differently. Specifically, the MFVB iterations were terminated when the relative change in $\log p(\mathbf{y}; q)$ fell below 10^{-7} , whilst MCMC was based on a kept sample of 1000. In addition, the speed achieved by MFVB is traded off against slight accuracy losses which are invoked by the product restriction given in (3.8).

However, even when allowing for these slight differences, the results clearly demonstrate that MFVB is at least approximately 200 times faster than MCMC across the models. Therefore, we conclude that a model which takes minutes to run using MCMC, will only take seconds to run using our MFVB algorithm.

	MCMC	MFVB
setting A	(1225.49, 1228.45)	(1.69, 1.87)
setting B	(1232.96, 1238.53)	(4.95, 5.89)
setting C	(1185.45, 1188.01)	(3.29, 4.67)
setting D	(1217.77, 1364.56)	(1.26, 1.38)

Table 3.3: 99% Wilcoxon confidence intervals based on run times in seconds for MCMC and MFVB fitting.

3.5 Application

We have seen the time benefits and good accuracy results for MFVB from the results discussed thus far using simulated data. We next evaluate the performance of MFVB using real data in order to understand its potential application.

Firstly, we consider two nonparametric regression examples without taking heteroscedasticity into account. Secondly, we use our new MFVB methodology by applying Algorithm 2 to the same regression examples. This enables us to clearly understand the implications of ignoring heteroscedasticity, as it might on occasion be tempting to immediately fit the model in a naïve fashion. Our process involved fitting a cubic smoothing spline to two nonparametric regression examples, based on the R function `smooth.spline()` (R Core

3.5. APPLICATION

Team, 2013) with default settings. Figure 3.6 illustrates the fits. For the interested reader, a complete description of the two data sets that were used can be found in Sections 2.7 and 5.3 of Ruppert *et al.* (2003).

As can be evidenced from the standardized residual plots in Figure 3.6, there exists significant heteroscedasticity. Ignoring the heteroscedastic nature of these data sets will lead to erroneous inferential statements, such as prediction intervals. In order to account

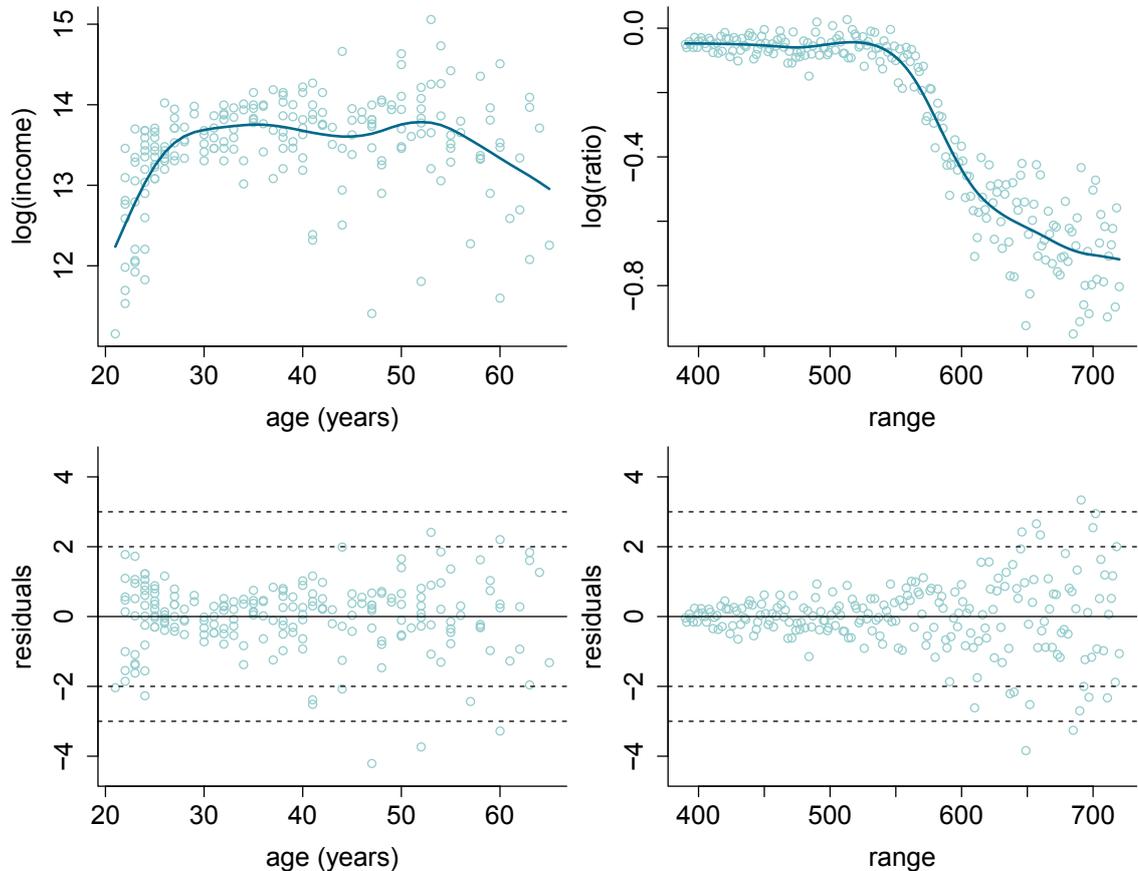


Figure 3.6: *Two example nonparametric regression fits and corresponding standardized residual plots. Horizontal dashed lines at ± 2 and ± 3 aid assessment of standard normality of the standardized residuals.*

for the significant heteroscedasticity as suggested by Figure 3.6, we also fitted (3.5) using Algorithm 2 to both these data sets. This is shown in Figure 3.7. The top panels illustrate the fitted mean functions with pointwise 95% credible sets. The middle panels illustrate the fitted standard deviation functions, corresponding to \sqrt{g} in the notation of (3.2), and corresponding 95% credible sets. The model was fitted using MFVB corresponding to Algorithm 2 and MCMC based on BUGS (?) with a burnin of 5000, a kept sample of

5000 and a thinning factor of 5. The MFVB iterations were terminated when the relative change in $\log \underline{p}(\mathbf{y}; q)$ fell below 10^{-7} . This value is determined by graphically checking if convergence is achieved and can change depending on the type of application. We also conducted a comprehensive check to confirm that the relative change in the approximate marginal log-likelihood does not lead to early stopping. Data was generated from different scenarios, over several hundred runs and Bayes estimates of f and g were recorded at the stopping point and again with iterations continuing 25% beyond the stopping point. The differences in the estimates were negligible.

The agreement between the MFVB and the MCMC fits for the mean function is excellent. For the standard deviation function fit, the agreement between MFVB and MCMC is good, rather than excellent. Finally, the \sqrt{g} -standardized residual plots in the bottom panels of Figure 3.7, indicate proper accounting for the heteroscedasticity and agreement with normality.

3.6 Discussion

At the outset of this chapter, we sought to develop a variational Bayesian algorithm for semiparametric regression models that have a heteroscedastic component. We sought to do this because ignoring heteroscedasticity, especially amongst complex statistical models, could lead to problematic inferential results and prediction intervals. As a result of the work carried out in this chapter, we have successfully developed a variational Bayes methodology which incorporates a new modification of MFVB, known as *non-conjugate variational message passing*, involving a closed form algorithm for univariate heteroscedastic semiparametric regression. The methodology also applies to larger more complex models, as provided for by the locality property of mean field variational inference methods.

Comparisons between variational inference and the MCMC benchmark for both simulated and actual data, shows that this new methodology performs very well. In addition, we have been able to achieve significant time savings by using our new methodology when compared to MCMC. For instance, as was shown in Section 3.4.3, we witness at least a 200 fold increase in estimation speed. This is equivalent to saying that a model which takes minutes to run using MCMC, will only take seconds to run using our MFVB methodology.

Having considered the impact of heteroscedasticity and remedies via our new MFVB-type methodology, we next turn our attention to extensions of this work which we discuss

3.6. DISCUSSION

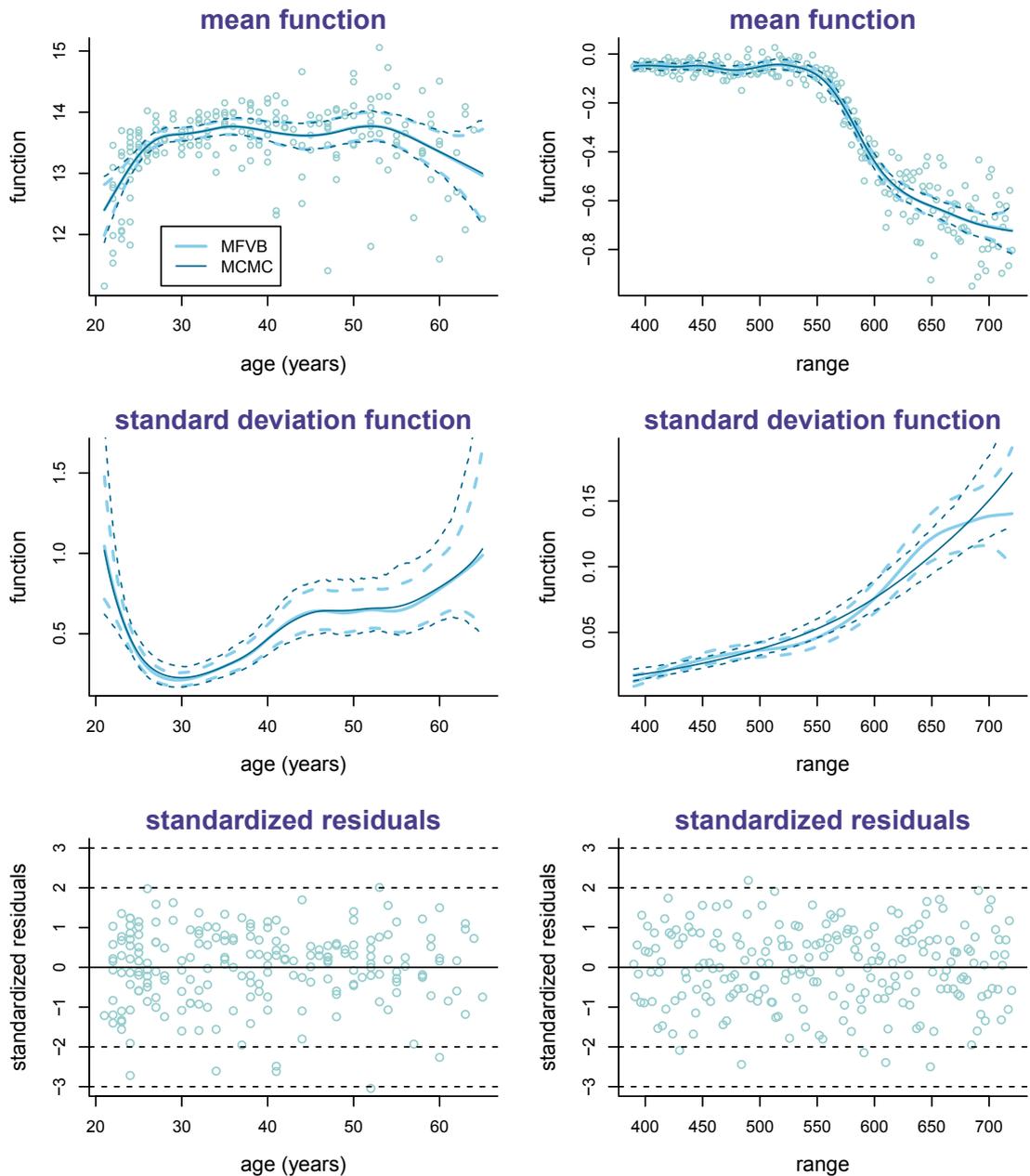


Figure 3.7: *Top panels: fitted mean functions for two example data sets from Wand & Jones (1995) and Ruppert et al. (2003). The solid curves are approximate pointwise posterior means whilst the dashed curves are corresponding pointwise 95% credible sets. The approximate fits are based on MFVB via Algorithm 2 and MCMC via BUGS. Middle panels: similar to top panels but for the standard deviation function. Bottom panels: standardized residual plots based on $\{y - \hat{f}(x_i)\} / \sqrt{\{\hat{g}(x_i)\}}$ where \hat{f} and \hat{g} are the MFVB-approximate Bayes estimates of f and g .*

in Chapter 4.

3.A Derivation of algorithm 2

Expressions for $\boldsymbol{\mu}_{q(\boldsymbol{\nu})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}$

First note that

$$\begin{aligned} p(\boldsymbol{\nu}|\text{rest}) &\propto p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega})p(\boldsymbol{\beta}) \\ &= \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})\right\} \times \exp\left\{-\frac{1}{2\sigma_u^2}\|\mathbf{u}\|^2\right\} \\ &\quad \times \exp\left\{-\frac{1}{2\sigma_\beta^2}\|\boldsymbol{\beta}\|^2\right\} + \text{const}, \end{aligned}$$

where ‘const’ denotes terms not depending on the argument of $\boldsymbol{\nu}$. Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\boldsymbol{\nu}|\text{rest}) &\propto -\frac{1}{2}(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu}) - \frac{1}{2\sigma_u^2}\|\mathbf{u}\|^2 - \frac{1}{2\sigma_\beta^2}\|\boldsymbol{\beta}\|^2 \\ &\propto -\frac{1}{2}(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})(\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu}) - \frac{1}{2}\boldsymbol{\nu}^\top \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \sigma_u^{-2}b\mathbf{I}_{K_u} \end{bmatrix} \boldsymbol{\nu} \\ &\propto -\frac{1}{2}\left\{\boldsymbol{\nu}^\top \mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})\mathbf{C}_\nu \boldsymbol{\nu} - 2\mathbf{C}_\nu^\top \boldsymbol{\nu}^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})\mathbf{y} \right. \\ &\quad \left. + \boldsymbol{\nu}^\top \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \sigma_u^{-2}b\mathbf{I}_{K_u} \end{bmatrix} \boldsymbol{\nu}\right\} \\ &= -\frac{1}{2}\left[\boldsymbol{\nu}^\top \left\{\mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})\mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \sigma_u^{-2}b\mathbf{I}_{K_u} \end{bmatrix}\right\} \boldsymbol{\nu} \right. \\ &\quad \left. - 2\boldsymbol{\nu}^\top \left\{\mathbf{C}_\nu^\top \text{diag}(e^{-\mathbf{C}_\omega \boldsymbol{\omega}})\mathbf{y}\right\}\right] + \text{const}. \end{aligned}$$

Now taking expectations with respect to all parameters except $\boldsymbol{\nu}$:

$$\begin{aligned} \log q^*(\boldsymbol{\nu}) &= -\frac{1}{2}\left[\boldsymbol{\nu}^\top \left\{\mathbf{C}_\nu^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})})\mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)}\mathbf{I}_{K_u} \end{bmatrix}\right\} \boldsymbol{\nu} \right. \\ &\quad \left. - 2\boldsymbol{\nu}^\top \left\{\mathbf{C}_\nu^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})})\mathbf{y}\right\}\right] + \text{const}, \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\psi}_{q(\boldsymbol{\omega})} &= E_q\{\exp(-\mathbf{C}_\omega \boldsymbol{\omega})\} \\ &= \exp\left\{-\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2}\text{diagonal}(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}_\omega^\top)\right\}. \end{aligned}$$

Completing the square gives

$$\log q^*(\boldsymbol{\nu}) = -\frac{1}{2}\left\{\boldsymbol{\nu} - \mathbf{C}_\nu^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})})\mathbf{y}\boldsymbol{\Omega}\right\}^\top \boldsymbol{\Omega} \left\{\boldsymbol{\nu} - \mathbf{C}_\nu^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})})\mathbf{y}\boldsymbol{\Omega}\right\}$$

where

$$\boldsymbol{\Omega} = \left(\mathbf{C}_\nu^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})})\mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2}\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)}\mathbf{I}_{K_u} \end{bmatrix}\right)^{-1}.$$

3.A. DERIVATION OF ALGORITHM 2

Therefore,

$q^*(\boldsymbol{\nu})$ is the $N(\boldsymbol{\mu}_{q(\boldsymbol{\nu})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\nu})})$ density function,

where

$$\boldsymbol{\mu}_{q(\boldsymbol{\nu})} = \boldsymbol{\Sigma}_{q(\boldsymbol{\nu})} \mathbf{C}_\nu^T \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \mathbf{y},$$

and

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})} = \left(\mathbf{C}_\nu^T \text{diag}\{\boldsymbol{\psi}_{q(\boldsymbol{\omega})}\} \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)} \mathbf{I}_{K_u} \end{bmatrix} \right)^{-1}.$$

Expressions for $B_{q(\sigma_u^2)}$ and $\mu_{q(1/\sigma_u^2)}$

$$\begin{aligned} p(\sigma_u^2 | \text{rest}) &\propto p(\mathbf{u} | \sigma_u^2) p(\sigma_u^2 | a_u) \\ &= |\sigma_u^2 \mathbf{I}_{K_u}|^{-1/2} \exp\left\{-\frac{1}{2\sigma_u^2} \|\mathbf{u}\|^2\right\} (\sigma_u^2)^{-\frac{1}{2}-1} \exp\left\{(1/a_u)/\sigma_u^2\right\} + \text{const.} \end{aligned}$$

Taking the logarithm of both sides, we have

$$\begin{aligned} \log p(\sigma_u^2 | \text{rest}) &= -\frac{1}{2} K_u \log(\sigma_u^2) - \frac{1}{2\sigma_u^2} \|\mathbf{u}\|^2 - \frac{3}{2} \log(\sigma_u^2) - (1/a_u)/\sigma_u^2 + \text{const.} \\ &= -\left\{\frac{1}{2}(K_u + 1) + 1\right\} \log(\sigma_u^2) - \left(\frac{1}{2} \|\mathbf{u}\|^2 + a_u^{-1}\right) / \sigma_u^2 + \text{const.} \end{aligned}$$

Taking expectations, we get

$$\begin{aligned} \log q^*(\sigma_u^2) &= E_q \left\{ \log p(\sigma_u^2 | \text{rest}) \right\} + \text{const.} \\ &= -\left\{\frac{1}{2}(K_u + 1) + 1\right\} \log(\sigma_u^2) - \left(\frac{1}{2} E_q \|\mathbf{u}\|^2 + \mu_{q(1/a_u)}\right) / \sigma_u^2 + \text{const} \end{aligned}$$

Noting that $E_q \|\mathbf{u}\|^2 = \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})})$, we have

$$\begin{aligned} \log q^*(\sigma_u^2) &= -\left\{\frac{1}{2}(K_u + 1) + 1\right\} \log(\sigma_u^2) - \left(\frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \right\} + \mu_{q(1/a_u)}\right) / \sigma_u^2 \\ &\quad + \text{const} \end{aligned}$$

Therefore,

$$q^*(\sigma_u^2) \propto (\sigma_u^2)^{-\left\{\frac{1}{2}(K_u+1)+1\right\}} \exp\left(-\frac{1}{\sigma_u^2} \left[\frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \right\} + \mu_{q(1/a_u)}\right]\right)$$

which is of the form of an Inverse-Gamma distribution with parameters

$$A_{q(\sigma_u^2)} = \frac{1}{2}(K_u + 1), \quad \text{and}$$

$$B_{q(\sigma_u^2)} = \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \right\} + \mu_{q(1/a_u)}.$$

3.A. DERIVATION OF ALGORITHM 2

In addition, using Result 1.4.3,

$$\mu_{q(1/\sigma_u^2)} = \frac{1}{2}(K_u + 1)/B_{q(\sigma_u^2)}.$$

Expressions for $B_{q(\sigma_v^2)}$ and $\mu_{q(1/\sigma_v^2)}$

The derivation of $B_{q(\sigma_v^2)}$ and $\mu_{q(1/\sigma_v^2)}$ is similar to that for $q^*(\sigma_u^2)$. The optimal q density satisfies

$$q^*(\sigma_v^2) \propto (\sigma_v^2)^{-\left\{\frac{1}{2}(K_v+1)+1\right\}} \exp\left(-\frac{1}{\sigma_v^2} \left[\frac{1}{2} \left\{\|\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})})\right\} + \mu_{q(1/a_v)}\right]\right)$$

which is of the form of an Inverse-Gamma distribution with parameters

$$A_{q(\sigma_v^2)} = \frac{1}{2}(K_v + 1) \quad \text{and}$$

$$B_{q(\sigma_v^2)} = \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})}) \right\} + \mu_{q(1/a_v)}.$$

Also,

$$\mu_{q(1/\sigma_v^2)} = \frac{1}{2}(K_v + 1)/B_{q(\sigma_v^2)}.$$

Expressions for $B_{q(a_u)}$ and $\mu_{q(1/a_u)}$

$$\begin{aligned} p(a_u|\text{rest}) &\propto p(\sigma_u^2|a_u)p(a_u) \\ &= (1/a_u)^{1/2} \exp\left\{- (1/a_u)/\sigma_u^2\right\} (a_u)^{-\frac{1}{2}-1} \exp\left\{- (1/A_u^2)/a_u\right\} + \text{const.} \end{aligned}$$

Taking the logarithm, we get

$$\log p(a_u|\text{rest}) = -2 \log(a_u) - (\sigma_u^{-2} + A_u^{-2})/a_u + \text{const.}$$

Taking expectations:

$$\begin{aligned} \log q^*(a_u) &= E_q \{ \log p(a_u|\text{rest}) \} + \text{const} \\ &= -2 \log(a_u) - E_q (\sigma_u^{-2} + A_u^{-2})/a_u + \text{const.} \\ &= -2 \log(a_u) - (\mu_{q(1/\sigma_u^2)} + A_u^{-2})/a_u + \text{const} \end{aligned}$$

Therefore,

$$q^*(a_u) \propto (a_u)^{-2} \exp\left\{- (\mu_{q(1/\sigma_u^2)} + A_u^{-2})/a_u\right\}$$

3.A. DERIVATION OF ALGORITHM 2

which is of the form of an Inverse-Gamma distribution with parameters

$$A_{q(a_u)} = 1, \quad \text{and}$$

$$B_{q(a_u)} = \mu_{q(1/\sigma_u^2)} + A_u^{-2}.$$

In addition, using Result 1.4.3,

$$\mu_{q(1/a_u)} = 1/B_{q(a_u)}.$$

Expressions for $B_{q(a_v)}$, and $\mu_{q(1/a_v)}$

The derivation of $B_{q(a_v)}$, and $\mu_{q(1/a_v)}$ is equivalent to that for $q^*(a_u)$. The optimal q density satisfies

$$q^*(a_v) \propto (a_v)^{-2} \exp \left\{ - \left(\mu_{q(1/\sigma_v^2)} + A_v^{-2} \right) / a_v \right\}$$

which is of the form of an Inverse-Gamma distribution with parameters

$$A_{q(a_v)} = 1, \quad \text{and}$$

$$B_{q(a_v)} = \mu_{q(1/\sigma_v^2)} + A_v^{-2}.$$

Also,

$$\mu_{q(a_v)} = \mu_{q(1/a_v)} = 1/B_{q(a_v)}.$$

Derivation of the $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ updates

The updates for $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ are based on maximisation of the current value of the marginal log-likelihood lower bound $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$ over these parameters using fixed point iteration. Wand (2014) shows that the updates reduce to

$$\begin{aligned} \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} &\leftarrow \left\{ -2 \text{vec}^{-1} \left(\left(\text{D}_{\text{vec}}(\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) S \right)^T \right) \right\}^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\omega})} &\leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \left(\text{D}_{\boldsymbol{\mu}_{q(\boldsymbol{\omega})}} S \right)^T. \end{aligned}$$

where

$$S \equiv E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\nu}, \boldsymbol{\omega}) + \log p(\boldsymbol{\omega} | \sigma_v^2) \right\},$$

D denotes derivative vector, as defined in Magnus & Neudecker (1999) and vec and vec^{-1} are as defined in Wand (2014). However, an equivalent expression for the first of these

3.A. DERIVATION OF ALGORITHM 2

updates is

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \leftarrow \left(-\mathbf{H}_{\boldsymbol{\mu}_{q(\boldsymbol{\omega})}} S \right)^{-1} \quad (3.14)$$

where \mathbf{H} denotes the Hessian matrix as defined in Magnus & Neudecker (1999). We work with this alternative form here. Proposition 3.A.1 gives justification for expression (3.14). Full details of the proof of Proposition 3.A.1 are given in Appendix 3.C.

Proposition 3.A.1. *Consider the $d \times 1$ random vector \mathbf{x} which follows the multivariate normal distribution and has density function of the form $f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$, then*

$$\frac{1}{2} \mathbf{H}_{\boldsymbol{\mu}} f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \text{vec}^{-1} \left(D_{\text{vec}(\boldsymbol{\Sigma})} f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \right).$$

An explicit expression for S is

$$\begin{aligned} S &\equiv -\mathbf{1}^T \mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} - \boldsymbol{\mu}_{q(r_v^2)}^T \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}_\omega^T \right) \right\} \\ &\quad - \frac{1}{2} \text{tr} \left(\begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \left(\boldsymbol{\mu}_{q(\boldsymbol{\omega})} \boldsymbol{\mu}_{q(\boldsymbol{\omega})}^T + \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \right) \right) \\ &\quad - \left(n + \frac{K}{2} + 1 \right) \log(2\pi) - \frac{K}{2} \log(\sigma_v^2) - \log(\sigma_\gamma^2). \end{aligned}$$

This leads to

$$\begin{aligned} d_{\boldsymbol{\mu}_{q(\boldsymbol{\omega})}} S &= -\mathbf{1}^T \mathbf{C}_\omega d\boldsymbol{\mu}_{q(\boldsymbol{\omega})} \\ &\quad + \boldsymbol{\mu}_{q(r_v^2)}^T \text{diag} \left(\exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}_\omega^T \right) \right\} \right) \mathbf{C}_\omega d\boldsymbol{\mu}_{q(\boldsymbol{\omega})} \\ &\quad - \boldsymbol{\mu}_{q(\boldsymbol{\omega})}^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} d\boldsymbol{\mu}_{q(\boldsymbol{\omega})} \\ &= \left(\left[\boldsymbol{\mu}_{q(r_v^2)} \odot \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}_\omega^T \right) \right\} - \mathbf{1} \right]^T \mathbf{C}_\omega \right. \\ &\quad \left. - \boldsymbol{\mu}_{q(\boldsymbol{\omega})}^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right) d\boldsymbol{\mu}_{q(\boldsymbol{\omega})}. \end{aligned}$$

3.A. DERIVATION OF ALGORITHM 2

Then, by Theorem 6, Chapter 5, of Magnus & Neudecker (1999),

$$\begin{aligned} \left(D_{\boldsymbol{\mu}_{q(\omega)}} S \right)^T &= \mathbf{C}_\omega^T \left[\boldsymbol{\mu}_{q(r_v^2)} \odot \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} - \mathbf{1} \right] \\ &\quad - \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \boldsymbol{\mu}_{q(\omega)}. \end{aligned}$$

Next,

$$\begin{aligned} d^2 \boldsymbol{\mu}_{q(\omega)} S &= \left(\boldsymbol{\mu}_{q(r_v^2)} \odot \exp \left[\left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} \right]^T \mathbf{C}_\omega \right. \\ &\quad \left. - \left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right) d\boldsymbol{\mu}_{q(\omega)} \\ &= \left[\left\{ \boldsymbol{\mu}_{q(r_v^2)} \odot \left(\text{diag} \left[\exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} \right] \right. \right. \right. \\ &\quad \left. \left. \left. \times \left(\mathbf{C}_\omega d\boldsymbol{\mu}_{q(\omega)} \right) \right) \right\}^T \mathbf{C}_\omega - \left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right] d\boldsymbol{\mu}_{q(\omega)} \end{aligned}$$

Then, using Result 1.4.11,

$$\begin{aligned} d^2 \boldsymbol{\mu}_{q(\omega)} S &= \left\{ \left(\text{diag} \left(\boldsymbol{\mu}_{q(r_v^2)} \right) \text{diag} \left[\exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} \right] \right. \right. \\ &\quad \left. \left. \times \left(-\mathbf{C}_\omega d\boldsymbol{\mu}_{q(\omega)} \right) \right)^T \mathbf{C}_\omega - \left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right\} d\boldsymbol{\mu}_{q(\omega)} \\ &= \left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \left(-\mathbf{C}_\omega^T \text{diag} \left[\boldsymbol{\mu}_{q(r_v^2)} \odot \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} \right. \right. \right. \\ &\quad \left. \left. \left. + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} \right] \right) \mathbf{C}_\omega d\boldsymbol{\mu}_{q(\omega)} \\ &\quad - \left(d\boldsymbol{\mu}_{q(\omega)} \right)^T \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} d\boldsymbol{\mu}_{q(\omega)} \end{aligned}$$

Therefore, by Theorem 6, Chapter 6, of Magnus & Neudecker (1999),

$$\begin{aligned} -\mathbf{H}_{\boldsymbol{\mu}_{q(\omega)}} S &= \mathbf{C}_\omega^T \text{diag} \left[\boldsymbol{\mu}_{q(r_v^2)} \odot \exp \left\{ \mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^T \right) \right\} \right] \mathbf{C}_\omega \\ &\quad + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix}. \end{aligned}$$

Combining these expressions leads to the updates in Algorithm 2.

3.B Derivation for the marginal log-likelihood lower bound

The expression for the lower bound on the marginal log-likelihood given in (3.14) is

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_q(\boldsymbol{\omega}), \boldsymbol{\Sigma}_q(\boldsymbol{\omega})) &= \frac{1}{2}(K_u + K_v + 4) - \frac{n}{2} \log(2\pi) + \log \Gamma\left(\frac{1}{2}(K_u + 1)\right) - 2 \log(\pi) \\
&+ \log \Gamma\left(\frac{1}{2}(K_v + 1)\right) - \log(A_u) - \log(A_v) - \frac{1}{2} \mathbf{1}^T (\mathbf{C}_\omega \boldsymbol{\mu}_q(\boldsymbol{\omega})) \\
&- \frac{1}{2} \mathbf{1}^T \left\{ \boldsymbol{\mu}_q(\sigma_u^2) \odot \exp(\boldsymbol{\psi}_q(\boldsymbol{\omega})) \right\} - \log(\sigma_\beta^2) - \log(\sigma_\gamma^2) \\
&+ \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\nu})| + \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\omega})| - \frac{1}{2\sigma_\beta^2} (\|\boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\beta}))) \\
&- \frac{1}{2\sigma_\gamma^2} (\|\boldsymbol{\mu}_q(\boldsymbol{\gamma})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\gamma}))) - \frac{1}{2}(K_u + 1) \log(B_{q(\sigma_u^2)}) \\
&- \frac{1}{2}(K_v + 1) \log(B_{q(\sigma_v^2)}) - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \\
&- \log(\mu_{q(1/\sigma_v^2)} + A_v^{-2}) + \mu_{q(1/\sigma_u^2)} \mu_{q(1/a_u)} + \mu_{q(1/\sigma_v^2)} \mu_{q(1/a_v)}.
\end{aligned}$$

The derivation of the logarithm lower bound on the marginal likelihood is due to

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\gamma}, \mathbf{v}, \sigma_u^2, \sigma_v^2, a_u, a_v) - \log q^*(\boldsymbol{\beta}, \mathbf{u}, \boldsymbol{\gamma}, \mathbf{v}, \sigma_u^2, \sigma_v^2, a_u, a_v) \right\} \\
&= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) - \log q^*(\boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) \right\}.
\end{aligned}$$

We can further factorise the true joint density since the distribution of \mathbf{y} does not depend on $\sigma_u^2, \sigma_v^2, a_u$ or a_v :

$$p(\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) = p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) p(\boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v).$$

Continuing in this fashion, we get

$$\begin{aligned}
p(\mathbf{y}, \boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) &= p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) p(\boldsymbol{\nu}|\sigma_u^2) p(\boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) \\
&= p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) p(\boldsymbol{\nu}|\sigma_u^2) p(\boldsymbol{\omega}|\sigma_v^2) p(\sigma_u^2, \sigma_v^2, a_u, a_v) \\
&= p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) p(\boldsymbol{\nu}|\sigma_u^2) p(\boldsymbol{\omega}|\sigma_v^2) p(\sigma_u^2) p(\sigma_v^2) p(a_u) p(a_v)
\end{aligned}$$

since the distribution of $\boldsymbol{\nu}$ does not depend on $\boldsymbol{\omega}, \sigma_v^2, a_u$ or a_v , and the distribution of $\boldsymbol{\omega}$ does not depend on a_u or a_v . Also, $\sigma_u^2, \sigma_v^2, a_u$ and a_v are independent.

The approximate joint density also factorises according to (3.10):

$$q(\boldsymbol{\nu}, \boldsymbol{\omega}, \sigma_u^2, \sigma_v^2, a_u, a_v) = q(\boldsymbol{\nu}) q(\boldsymbol{\omega}) q(\sigma_u^2) q(\sigma_v^2) q(a_u) q(a_v).$$

Using the above factorisations, the new expression for the logarithm lower bound is

3.B. DERIVATION FOR THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

now:

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= E_q \{ \log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) \} + E_q \{ \log p(\boldsymbol{\nu}|\sigma_u^2) - q^*(\boldsymbol{\nu}) \} \\ &\quad + E_q \{ \log p(\boldsymbol{\omega}|\sigma_v^2) - \log q^*(\boldsymbol{\omega}) \} + E_q \{ \log p(\sigma_u^2|a_u) - \log q^*(\sigma_u^2) \} \\ &\quad + E_q \{ \log p(\sigma_v^2|a_v) - \log q^*(\sigma_v^2) \} + E_q \{ \log p(a_u) - \log q^*(a_u) \} \\ &\quad + E_q \{ \log p(a_v) - \log q^*(a_v) \}. \end{aligned}$$

The derivation of each component is given as follows.

Expression for $E_q \{ \log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) \}$

Firstly,

$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) &= (2\pi)^{-n/2} |\text{diag} \{ \exp(\mathbf{C}_\omega \boldsymbol{\omega}) \}|^{-1/2} \\ &\quad \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag} \{ \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \} \right. \\ &\quad \left. (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu}) \right\}. \end{aligned}$$

Note that

$$\begin{aligned} |\text{diag} \{ \exp(\mathbf{C}_\omega \boldsymbol{\omega}) \}| &= \begin{bmatrix} \exp \{ (\mathbf{C}_\omega \boldsymbol{\omega})_1 \} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \exp \{ (\mathbf{C}_\omega \boldsymbol{\omega})_n \} \end{bmatrix} \\ &= \exp \{ (\mathbf{C}_\omega \boldsymbol{\omega})_1 \} \times \dots \times \exp \{ (\mathbf{C}_\omega \boldsymbol{\omega})_n \} \\ &= \exp \{ \mathbf{1}^\top (\mathbf{C}_\omega \boldsymbol{\omega}) \} \end{aligned}$$

and

$$\begin{aligned} &-\frac{1}{2} (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag} \{ \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \} (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\nu}) \\ &= -\frac{1}{2} (\mathbf{C}_\nu \boldsymbol{\nu})^\top \text{diag} \{ \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \} (\mathbf{C}_\nu \boldsymbol{\nu}) \\ &= -\frac{1}{2} r_\nu^\top \text{diag} \{ \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \} r_\nu \\ &= -\frac{1}{2} \mathbf{1}^\top \{ r_\nu^2 \odot \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \} \end{aligned}$$

where

$$\begin{aligned} r_\nu &= \mathbf{C}_\nu \boldsymbol{\nu}, \quad \text{and} \\ r_\nu^2 &= \text{diagonal} \{ (\mathbf{C}_\nu \boldsymbol{\nu}) (\mathbf{C}_\nu \boldsymbol{\nu})^\top \}. \end{aligned}$$

Therefore

$$\log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \mathbf{1}^\top (\mathbf{C}_\omega \boldsymbol{\omega}) - \frac{1}{2} \mathbf{1}^\top \{ r_\nu^2 \odot \exp(-\mathbf{C}_\omega \boldsymbol{\omega}) \}.$$

Taking expectations we get

$$E_q \{ \log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\omega}) \} = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \mathbf{1}^\top (\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})}) - \frac{1}{2} \mathbf{1}^\top \{ \boldsymbol{\mu}_{q(r_\nu^2)} \odot \exp(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \},$$

where

$$\begin{aligned}\boldsymbol{\mu}_{q(\tau_v^2)} &= E_q \left[\text{diagonal} \left\{ (\mathbf{C}_\nu \boldsymbol{\nu}) (\mathbf{C}_\nu \boldsymbol{\nu})^\top \right\} \right] \\ &= \text{diagonal} \left[E_q \left\{ (\mathbf{C}_\nu \boldsymbol{\nu}) (\mathbf{C}_\nu \boldsymbol{\nu})^\top \right\} \right] \\ &= \text{diagonal} \left\{ (\mathbf{C}_\nu \boldsymbol{\mu}_{q(\nu)}) (\mathbf{C}_\nu \boldsymbol{\mu}_{q(\nu)})^\top + \mathbf{C}_\nu \boldsymbol{\Sigma}_{q(\nu)} \mathbf{C}_\nu^\top \right\}.\end{aligned}$$

Expression for $E_q \{ \log p(\boldsymbol{\nu} | \sigma_u^2) - \log q^*(\boldsymbol{\nu}) \}$

First note that

$$\log p(\boldsymbol{\nu} | \sigma_u^2) = -\log(2\pi) - \left(\frac{k}{2} + 1\right) \log(2\pi) - \log(\sigma_\beta^2) - \frac{k}{2} \log(\sigma_u^2) - \frac{1}{2\sigma_\beta^2} \|\boldsymbol{\beta}\|^2 - \frac{1}{2\sigma_u^2} \|\mathbf{u}\|^2.$$

Taking expectations we get

$$\begin{aligned}E_q \{ \log p(\boldsymbol{\nu} | \sigma_u^2) \} &= -\log(2\pi) - \left(\frac{k}{2} + 1\right) \log(2\pi) - \log(\sigma_\beta^2) - \frac{k}{2} E_q \{ \log(\sigma_u^2) \} \\ &\quad - \frac{1}{2\sigma_\beta^2} E_q (\|\boldsymbol{\beta}\|^2) - \frac{1}{2\sigma_u^2} E_q (\|\mathbf{u}\|^2) \\ &= -\log(2\pi) - \left(\frac{k}{2} + 1\right) \log(2\pi) - \log(\sigma_\beta^2) - \frac{k}{2} E_q \{ \log(\sigma_u^2) \} \\ &\quad - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} - \frac{1}{2} \mu_{q(1/\sigma_u^2)} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \right\}.\end{aligned}$$

Also,

$$E_q \{ \log q^*(\boldsymbol{\nu}) \} = \frac{1}{2} (K_u + 2) \log(2\pi) + \frac{1}{2} (K_u + 2) + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\nu)}|.$$

Thus, we get

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

Expression for $E_q \{ \log p(\boldsymbol{\omega} | \sigma_v^2) - \log q^*(\boldsymbol{\omega}) \}$

Similar to the previous derivation we get

$$\begin{aligned}E_q \{ \log p(\boldsymbol{\omega} | \sigma_v^2) - \log q^*(\boldsymbol{\omega}) \} &= -\log(\sigma_\gamma^2) - \frac{1}{2} K_v E_q (\log(\sigma_v^2)) + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}| \\ &\quad - \frac{1}{2\sigma_\gamma^2} \left\{ \|\boldsymbol{\mu}_{q(\gamma)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\gamma)}) \right\} + \frac{1}{2} (K_v + 2) \\ &\quad - \frac{1}{2} \mu_{q(1/\sigma_v^2)} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})}) \right\}.\end{aligned}$$

Expression for $E_q \{ \log p(\sigma_u^2 | a_u) - \log q^*(\sigma_u^2) \}$

Firstly,

$$\begin{aligned} \log p(\sigma_u^2|a_u) - \log q^*(\sigma_u^2) &= \log \left\{ \frac{\left(\frac{1}{a_u}\right)^{\frac{1}{2}}}{\Gamma\left(\frac{1}{2}\right)} (\sigma_u^2)^{-\frac{1}{2}-1} \exp\left(-\frac{1}{\sigma_u^2}\right) \right\} \\ &\quad - \log \left\{ \frac{B_q(\sigma_u^2)^{\frac{1}{2}(K_u+1)}}{\Gamma\left(\frac{1}{2}(K_u+1)\right)} (\sigma_u^2)^{-\frac{1}{2}(K_u+1)-1} \exp\left(-\frac{B_q(\sigma_u^2)}{\sigma_u^2}\right) \right\} \\ &= -\frac{1}{2} \log(a_u) + \log \Gamma\left\{\frac{1}{2}(K_u+1)\right\} - \frac{1}{2} \log(\pi) + \frac{1}{2} K_u \log(\sigma_u^2) \\ &\quad + (1/\sigma_u^2) (B_q(\sigma_u^2) - a_u^{-1}) - \frac{1}{2} (K_u+1) \log(B_q(\sigma_u^2)). \end{aligned}$$

Taking expectations we get

$$\begin{aligned} E_q \{ \log p(\sigma_u^2|a_u) - \log q^*(\sigma_u^2) \} &= -\frac{1}{2} E_q \{ \log(a_u) \} + \log \Gamma\left\{\frac{1}{2}(K_u+1)\right\} - \frac{1}{2} \log(\pi) \\ &\quad + \mu_{q(1/\sigma_u^2)} (B_q(\sigma_u^2) - \mu_{q(1/a_u)}) + \frac{1}{2} K_u E_q \{ \log(\sigma_u^2) \} \\ &\quad - \frac{1}{2} (K_u+1) \log(B_q(\sigma_u^2)). \end{aligned}$$

Expression for $E_q \{ \log p(\sigma_v^2|a_v) - \log q^*(\sigma_v^2) \}$

Similar to the previous derivation we get

$$\begin{aligned} E_q \{ \log p(\sigma_v^2|a_v) - \log q^*(\sigma_v^2) \} &= -\frac{1}{2} E_q \{ \log(a_v) \} + \log \Gamma\left\{\frac{1}{2}(K_v+1)\right\} - \frac{1}{2} \log(\pi) \\ &\quad + \mu_{q(1/\sigma_v^2)} (B_q(\sigma_v^2) - \mu_{q(1/a_v)}) + \frac{1}{2} K_v E_q \{ \log(\sigma_v^2) \} \\ &\quad - \frac{1}{2} (K_v+1) \log(B_q(\sigma_v^2)). \end{aligned}$$

Expression $E_q \{ \log p(a_u) - \log q^*(a_u) \}$

The difference of the logarithms are

$$\begin{aligned} \log p(a_u) - \log q^*(a_u) &= \log \left\{ \frac{\left(\frac{1}{A_u}\right)^{\frac{1}{2}}}{\Gamma\left(\frac{1}{2}\right)} (a_u)^{-\frac{1}{2}-1} \exp\left(-\frac{1}{a_u}\right) \right\} \\ &\quad - \log \left\{ \frac{B_q(a_u)}{\Gamma(1)} (a_u)^{-2} \exp\left(-\frac{B_q(a_u)}{a_u}\right) \right\} \\ &= -\log(A_u) - \log \Gamma\left(\frac{1}{2}\right) + \frac{1}{2} \log(a_u) + (1/a_u) (B_q(a_u) - A_u^{-2}) \\ &\quad - \log(B_q(a_u)). \end{aligned}$$

Then, substituting $B_q(a_u)$ from Algorithm 2 gives

$$\begin{aligned} \log p(a_u) - \log q^*(a_u) &= -\log(A_u) - \frac{1}{2} \log(\pi) + \frac{1}{2} \log(a_u) + \mu_{q(1/\sigma_u^2)}/a_u \\ &\quad - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}). \end{aligned}$$

Taking expectations, we get

$$\begin{aligned} E_q \{ \log p(a_u) - \log q^*(a_u) \} &= -\log(A_u) - \frac{1}{2} \log(\pi) + \frac{1}{2} E_q \{ \log(a_u) \} + \mu_{q(1/\sigma_u^2)} \mu_{q(1/a_u)} \\ &\quad - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}). \end{aligned}$$

Expression $E_q \{ \log p(a_v) - \log q^*(a_v) \}$

Similar to the previous derivation we get

$$E_q \{ \log p(a_v) - \log q^*(a_v) \} = -\log(A_v) - \frac{1}{2} \log(\pi) + \frac{1}{2} E_q \{ \log(a_v) \} + \mu_{q(1/\sigma_v^2)} \mu_{q(1/a_v)} - \log(\mu_{q(1/\sigma_v^2)} + A_v^{-2}).$$

Adding all expressions together, we get the lower bound expression

$$\begin{aligned} \log p(\underline{\mathbf{y}}; q, \boldsymbol{\mu}_q(\boldsymbol{\omega}), \boldsymbol{\Sigma}_q(\boldsymbol{\omega})) &= \frac{1}{2} (K_u + K_v + 4) - \frac{n}{2} \log(2\pi) + \log \Gamma\left(\frac{1}{2}(K_u + 1)\right) - 2 \log(\pi) \\ &+ \log \Gamma\left(\frac{1}{2}(K_v + 1)\right) - \log(A_u) - \log(A_v) - \frac{1}{2} \mathbf{1}^T (\mathbf{C}_\omega \boldsymbol{\mu}_q(\boldsymbol{\omega})) \\ &- \frac{1}{2} \mathbf{1}^T \{ \boldsymbol{\mu}_{q(\sigma_u^2)} \odot \exp(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \} - \log(\sigma_\beta^2) - \log(\sigma_\gamma^2) \\ &+ \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\nu})}| + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}| - \frac{1}{2\sigma_\beta^2} (\| \boldsymbol{\mu}_{q(\boldsymbol{\beta})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})})) \\ &- \frac{1}{2\sigma_\gamma^2} (\| \boldsymbol{\mu}_{q(\boldsymbol{\gamma})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\gamma})})) - \frac{1}{2} (K_u + 1) \log(B_{q(\sigma_u^2)}) \\ &- \frac{1}{2} (K_v + 1) \log(B_{q(\sigma_v^2)}) - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \\ &- \log(\mu_{q(1/\sigma_v^2)} + A_v^{-2}) + \mu_{q(1/\sigma_u^2)} \mu_{q(1/a_u)} + \mu_{q(1/\sigma_v^2)} \mu_{q(1/a_v)}. \end{aligned}$$

Note the following cancellations:

$$-\frac{1}{2} \mu_{q(1/\sigma_u^2)} \{ \| \boldsymbol{\mu}_{q(\mathbf{u})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \} + \mu_{q(1/\sigma_u^2)} (B_{q(\sigma_u^2)} - \mu_{q(1/a_u)}),$$

where

$$B_{q(\sigma_u^2)} = \frac{1}{2} \{ \| \boldsymbol{\mu}_{q(\mathbf{u})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \} + \mu_{q(1/a_u)}$$

and

$$-\frac{1}{2} \mu_{q(1/\sigma_v^2)} \{ \| \boldsymbol{\mu}_{q(\mathbf{v})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})}) \} + \mu_{q(1/\sigma_v^2)} (B_{q(\sigma_v^2)} - \mu_{q(1/a_v)}),$$

where

$$B_{q(\sigma_v^2)} = \frac{1}{2} \{ \| \boldsymbol{\mu}_{q(\mathbf{v})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})}) \} + \mu_{q(1/a_v)}.$$

3.C Proof of Proposition 3.A.1

We begin by expanding the first differential with respect to $\text{vec}(\Sigma)$ in the right-hand side of Proposition 3.A.1:

$$\begin{aligned} d_{\text{vec}(\Sigma)}f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) &= \left(d|\Sigma|^{-\frac{1}{2}} \right) \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &\quad + |\Sigma|^{-\frac{1}{2}} d \left[\exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \right] \\ &= -\frac{1}{2} |\Sigma|^{-\frac{3}{2}} d|\Sigma| \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &\quad + |\Sigma|^{-\frac{1}{2}} \left(-\frac{1}{2} \text{tr} \left\{ (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top d\Sigma^{-1} \right\} \right) \\ &\quad \times \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Using Results 1.4.1 (b) and (c), we get

$$\begin{aligned} d_{\text{vec}(\Sigma)}f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) &= -\frac{1}{2} |\Sigma|^{-\frac{3}{2}} |\Sigma| \text{tr}(\Sigma^{-1} d\Sigma) \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &\quad - \frac{1}{2} |\Sigma|^{-\frac{1}{2}} \text{tr} \left\{ (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top (-\Sigma^{-1} (d\Sigma) \Sigma^{-1}) \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Next, using Results 1.4.13 and 1.4.1 (a), we see that

$$\begin{aligned} d_{\text{vec}(\Sigma)}f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) &= -\frac{1}{2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \text{vec}(\Sigma^{-1}) d\text{vec}(\Sigma) \\ &\quad + \frac{1}{2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &\quad \times \text{vec}(\Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}) d\text{vec}(\Sigma). \end{aligned}$$

Thus,

$$\begin{aligned} \text{vec}^{-1} \left(D_{\text{vec}(\Sigma)}f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \right) &= -\frac{1}{2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \Sigma^{-1} \\ &\quad + \frac{1}{2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \Sigma^{-1} \\ &\quad \times (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} \\ &= -\frac{1}{2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &\quad \times \Sigma^{-1} \left\{ \mathbf{I} - (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} \right\}. \end{aligned}$$

3.C. PROOF OF PROPOSITION 3.A.1

Now expanding the left-hand side of Proposition 3.A.1 we have

$$\begin{aligned}
d_{\boldsymbol{\mu}}f(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \left(d|\boldsymbol{\Sigma}|^{-\frac{1}{2}}\right) \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \\
&\quad + |\boldsymbol{\Sigma}|^{-\frac{1}{2}} d\left[\exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\}\right] \\
&= |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \left[-\frac{1}{2}\left\{-(d\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}) + (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(-d\boldsymbol{\mu})\right\}\right] \\
&\quad \times \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \\
&= |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} d\boldsymbol{\mu}.
\end{aligned}$$

Next, we take the second differential with respect to $\boldsymbol{\mu}$:

$$\begin{aligned}
d_{\boldsymbol{\mu}}^2 f(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= d\boldsymbol{\mu} \left[|\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \right] d\boldsymbol{\mu} \\
&= \left[d\boldsymbol{\mu} \left\{ |\boldsymbol{\Sigma}|^{-\frac{1}{2}} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \right\} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \right. \\
&\quad \left. + |\boldsymbol{\Sigma}|^{-\frac{1}{2}} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} d\boldsymbol{\mu} \left[\exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \right] \right] d\boldsymbol{\mu} \\
&= \left[|\boldsymbol{\Sigma}|^{-\frac{1}{2}} (-d\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \right. \\
&\quad \left. + |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (d\boldsymbol{\mu})^\top \right. \\
&\quad \left. \times \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \right] d\boldsymbol{\mu} \\
&= (d\boldsymbol{\mu})^\top \left[-|\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \boldsymbol{\Sigma}^{-1} \right. \\
&\quad \left. \times \left\{ \boldsymbol{I} - (\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \right\} \right] d\boldsymbol{\mu}.
\end{aligned}$$

Therefore,

$$\mathbf{H}_{\boldsymbol{\mu}}f(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -|\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right\} \boldsymbol{\Sigma}^{-1} \left\{ \boldsymbol{I} - (\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \right\}.$$

Proposition 3.A.1 follows immediately. This proposition was inspired by a result in the appendix of Opper & Archambeau (2009).

Chapter 4

Heteroscedastic semiparametric regression extensions

4.1 Introduction

The 2010s has engendered an unprecedented explosion of data, spurring a need for faster and increasingly flexible techniques to extrapolate and make sense of such data. The use of semiparametric regression (e.g. Ruppert *et al.*, 2003, 2009) provides a significant collection of tractable statistical techniques, however more inline with datasets comprising moderate sample sizes and batch processing. Of recent, Luts *et al.* (2013) developed algorithms to perform semiparametric regression analysis in real time, where the data is processed as it is collected and promptly made available through modern technologies. This new nonparametric regression methodology is especially geared toward high volume/velocity data. Part of this chapter extends their approach to adjust for possible heteroscedasticity.

Chapter 3 has shown that ignorance of significant heteroscedasticity is likely to lead to inaccurate inferential statistics. A remedy for this is the employment of heteroscedastic nonparametric regression. In this chapter we explore two other extensions to the univariate nonparametric regression model. Here we instead focus on bivariate predictor nonparametric regression and additive models. The modularity of our approach is utilised here since the extension to bivariate and additive models has been relatively straightforward.

The content of this chapter is published as: Menictas, M and Wand. M.P. (2015). Variational inference for heteroscedastic semiparametric regression. *Australian and New Zealand Journal of Statistics*, **57**, Number 1, 119–138. This research was also presented at the *Australian Statistical Conference in conjunction with the Institute of Mathematical Statistics Annual Meeting*, Sydney, 2014.

Much like in Chapter 3 we take a Bayesian approach here as well by implementing a non-conjugate MFVB strategy.

The following section describes the bivariate nonparametric regression extension methodology and illustrates this approach on simulated and real data. Accuracy scores and time comparisons are also given. Section 4.3 describes the additive model extension and also provides some illustrations on real data. Finally in Section 4.4 we extend the work of Luts *et al.* (2013) to accommodate heteroscedasticity.

4.2 Bivariate predictor nonparametric regression

Bivariate nonparametric regression is of particular interest to a number of disciplines, such as geology, oceanography, public health and mining. This diverse area of statistics deals with tractable smoothing of *point clouds* where we would be interested in the mean response as a bivariate function of, for instance, the longitude and latitude of a certain geographical location. Bivariate smoothing and *geostatistics* have many similarities, for instance, geostatistics is concerned with converting geographically referenced observations to increasingly more legible maps (e.g. Cressie, 1993) and hence deals with a bivariate predictor component.

As seen in previous chapters, the use of penalized spline smoothing requires a set of basis functions that enables the possibility of nonlinear structure. Bivariate smoothing on the other hand requires bivariate basis functions. This extension can be treated in at least two specific ways: (i) *tensor product bases* and (ii) *radial basis functions*. More detail of these approaches is given in Ruppert *et al.* (2003). A possible difficulty of *tensor product bases* is their dependence on the position of the coordinate axes, causing the results to change if the geographical locations were measured on axes with different positions. On the other hand, *radial basis functions* are invariant to rotation of the axes. In nongeographical applications rotational invariance is not an obstacle, but for geographical data it is important to know that the answers are not affected by axis position. Hence, for this extension we use *radial basis functions* for bivariate smoothing, in particular the low-rank thin plate spline bases.

4.2.1 Model description

Algorithm 2 is relatively easy to extend to its bivariate extension. The data are of the form

$$(\mathbf{x}_i, y_i), \quad \mathbf{x}_i \in \mathbb{R}^2, \quad y_i \in \mathbb{R}. \quad (4.1)$$

In classical geostatistics applications the \mathbf{x}_i s represent geographical locations, such as latitude and longitude. However, in (4.1) the \mathbf{x}_i s could also represent pairs of non-geographical measurements. The bivariate version of (3.2) is

$$y_i \stackrel{\text{ind.}}{\sim} N(f(\mathbf{x}_i), g(\mathbf{x}_i)), \quad 1 \leq i \leq n, \quad (4.2)$$

where f and g are real-valued smooth functions on \mathbb{R}^2 . The bivariate analogue of (3.3) is

$$\begin{aligned} f(\mathbf{x}) &= \beta_0 + \boldsymbol{\beta}_1^\top \mathbf{x} + \sum_{k=1}^{K_u} u_k z_k^u(\mathbf{x}), & u_k &\stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2) \\ \text{and } g(\mathbf{x}) &= \exp\left(\gamma_0 + \boldsymbol{\gamma}_1^\top \mathbf{x} + \sum_{k=1}^{K_v} v_k z_k^v(\mathbf{x})\right), & v_k &\stackrel{\text{ind.}}{\sim} N(0, \sigma_v^2), \end{aligned} \quad (4.3)$$

where $\boldsymbol{\beta}_1$ and $\boldsymbol{\gamma}_1$ are each 2×1 vectors of fixed effects. The functions $\{z_k^u : 1 \leq k \leq K_u\}$ and $\{z_k^v : 1 \leq k \leq K_v\}$ now represent bivariate spline basis functions. A reasonable default for the z_k^u (Ruppert *et al.*, 2003) is the low-rank thin plate spline basis with k th element:

$$z_k^u(\mathbf{x}) = r(\|\mathbf{x} - \boldsymbol{\kappa}_k^u\|) \left[r(\|\boldsymbol{\kappa}_k^u - \boldsymbol{\kappa}_{k'}^u\|) \right]_{1 \leq k, k' \leq K_u}^{-1/2}, \quad (4.4)$$

where $\boldsymbol{\kappa}_1^u, \dots, \boldsymbol{\kappa}_{K_u}^u$ is a set of bivariate knot locations that efficiently cover the space of the \mathbf{x}_i s and $r(x) \equiv x^2 \log(x)$. The default z_k^v s have an analogous definition.

Comparing this set-up to the univariate nonparametric heteroscedastic regression model, treated in Chapter 3, the only differences are the basis functions and their coefficients. Hence, Algorithm 2 can be used to fit the bivariate nonparametric heteroscedastic regression model by replacing $\boldsymbol{\nu}$, $\boldsymbol{\omega}$, \mathbf{C}_ν and \mathbf{C}_ω from Section 3.2 with

$$\boldsymbol{\nu} \equiv \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta}_1 \\ \mathbf{u} \end{bmatrix}, \quad \boldsymbol{\omega} \equiv \begin{bmatrix} \gamma_0 \\ \boldsymbol{\gamma}_1 \\ \mathbf{v} \end{bmatrix}, \quad \mathbf{C}_\nu \equiv \left[\begin{array}{c} 1 \quad \mathbf{x}_i^T \\ \mathbf{1}_{1 \leq k \leq K_u} \quad z_k^u(\mathbf{x}_i) \end{array} \right]_{1 \leq i \leq n} \quad \text{and} \quad \mathbf{C}_\omega \equiv \left[\begin{array}{c} 1 \quad \mathbf{x}_i^T \\ \mathbf{1}_{1 \leq k \leq K_v} \quad z_k^v(\mathbf{x}_i) \end{array} \right]_{1 \leq i \leq n}.$$

We have carried out extensive assessment of the bivariate version of Algorithm 2 to assess accuracy and computing time of model (4.2). Sections 4.2.2 and 4.2.3 give the details of a simulation study and real data example.

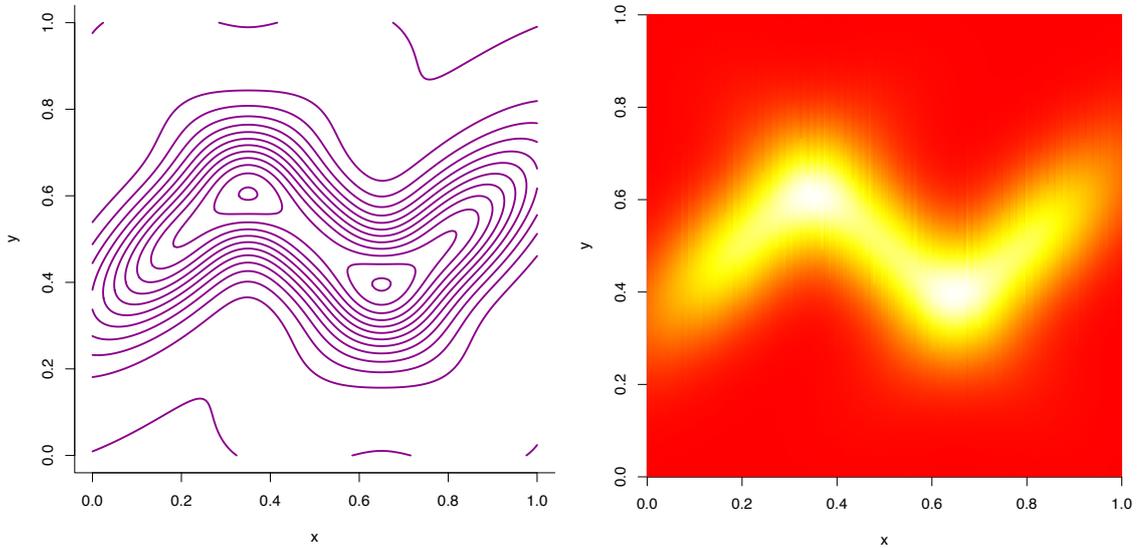


Figure 4.1: *Left panel: Contour plot illustrating the simulation setting given by (4.5). Right panel: Heat plot corresponding to (4.5).*

4.2.2 Simulation study

We consider the following true mean and log variance functions for our simulation setting. This is illustrated in Figure 4.1.

$$\begin{aligned}
 f(x_1, x_2) &= 0.199 + 12.924 \left\{ \phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} 0.2 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.360 & 0.252 \\ 0.252 & 0.360 \end{bmatrix} \right) \right. \\
 &\quad + \phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} 0.375 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.360 & 0.252 \\ 0.252 & 0.360 \end{bmatrix} \right) \\
 &\quad \left. + \phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.360 & -0.252 \\ -0.252 & 0.360 \end{bmatrix} \right) \right\}, \\
 \log g(x_1, x_2) &= -5.470 + 0.873 \cos(0.3\pi(x_1 + x_2)).
 \end{aligned} \tag{4.5}$$

From this setting we generated 1000 datasets of size $n = 200$. Each model corresponding to a new replication was fitted using MFVB corresponding to the bivariate version of Algorithm 2 and MCMC based on a burnin of 5000, kept sample of 5000 and thinning factor of 5. The MFVB iterations were terminated when the relative change in $\log p(\mathbf{y}; q, \boldsymbol{\mu}_q(\omega), \boldsymbol{\Sigma}_q(\omega))$ fell below 10^{-7} .

Figure 4.2 illustrates an accuracy assessment where the accuracy scores for each parameter of interest are summarized as a boxplot. The parameters on the horizontal axis of Figure 4.2 are the estimated approximate posterior density functions for f and g , evaluated at the sample quartiles of x_1 and x_2 . We use Q_1, Q_2 and Q_3 to denote these sample

4.2. BIVARIATE PREDICTOR NONPARAMETRIC REGRESSION

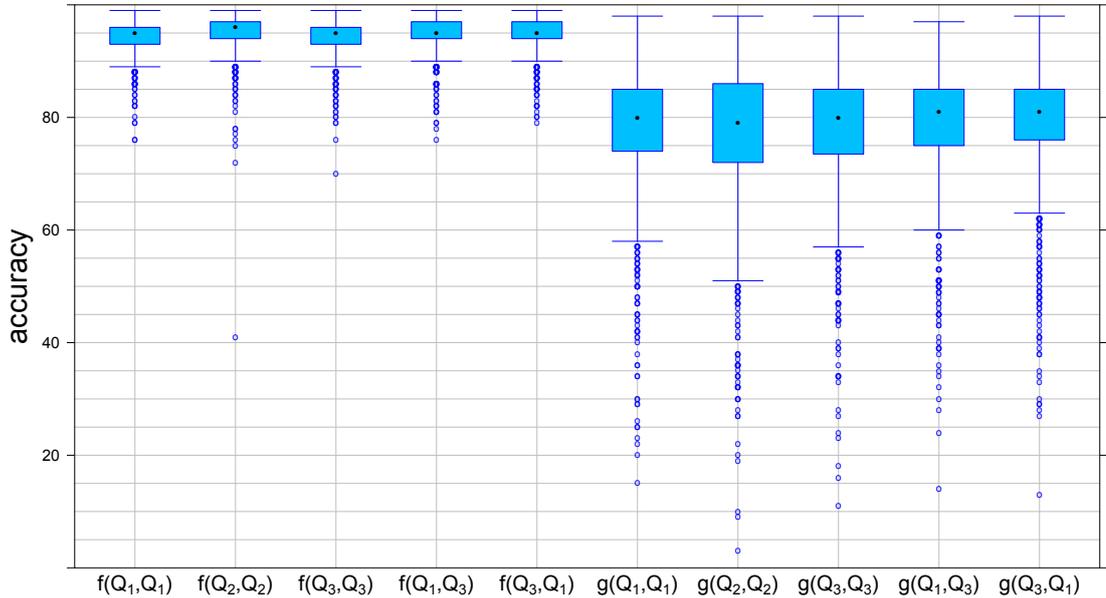


Figure 4.2: *Boxplot representation of accuracy values for MFVB against an MCMC benchmark, where the boxplots correspond to the mean function hexiles and variance function hexiles.*

quartiles. It is pleasing to see that most of the accuracy values for the bivariate mean function lie above 90% and above 70% for the bivariate variance function. Figure 4.3 shows visual assessment of the approximate posterior density functions as given by using MFVB in comparison with the MCMC benchmark for a single replication in the simulation study. MFVB accuracy is seen to be excellent for the mean function hexiles and quite good for the variance function hexiles.

The percentages of the true parameter coverage based on the approximate 95% credible sets attained from the MFVB posterior densities are given in Table 4.1. All of the mean function coverage values are showing excellent coverage of the truth, except for $f(Q_2, Q_2)$ and this may be due to the lack of data surrounding this region in our simulated datasets. Also, the variance function coverage values are doing well in terms of covering the truth and this is consistent with the accuracy scores achieved by MFVB. An assessment of the speed of the competing approaches has also been carried out. The run time of each replication was monitored for both MCMC and MFVB. These times are summarized in Table 4.2. This study was performed on a laptop computer (Intel Core i7 2.8GHz processor, 16 GBytes of random access memory). These results show that on average

4.2. BIVARIATE PREDICTOR NONPARAMETRIC REGRESSION

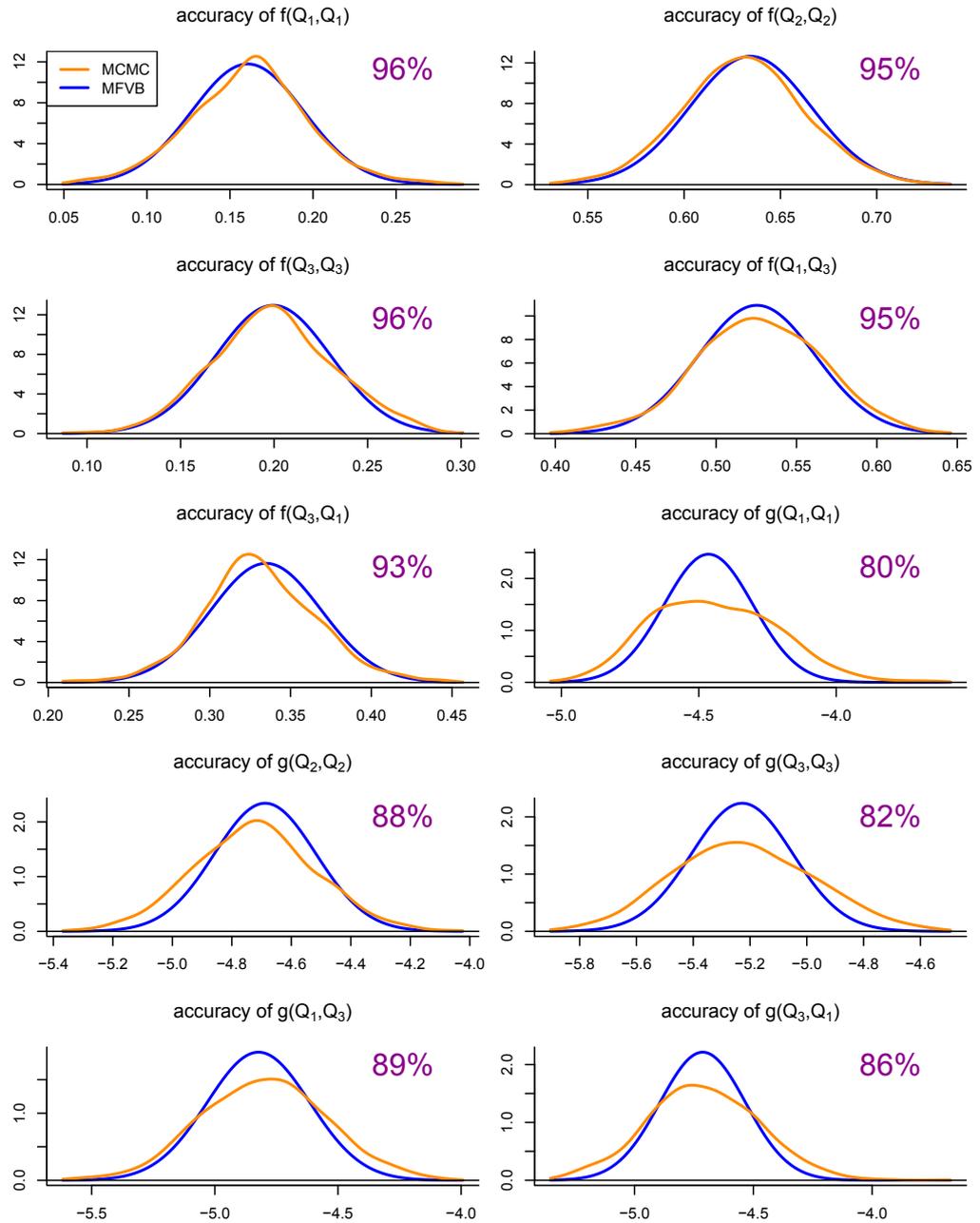


Figure 4.3: *Approximate posterior density functions obtained via MFVB and MCMC for a single replication of the simulation study. The pairs of density functions shown correspond to the mean and variance functions hexiles. The accuracy scores as defined in Section 1.9 are also displayed.*

4.2. BIVARIATE PREDICTOR NONPARAMETRIC REGRESSION

$f(Q_1, Q_1)$	$f(Q_2, Q_2)$	$f(Q_3, Q_3)$	$f(f(Q_1, Q_3))$	$f(Q_3, Q_1)$
94	68	92	90	91
$g(Q_1, Q_1)$	$g(Q_2, Q_2)$	$g(Q_3, Q_3)$	$g(Q_1, Q_3)$	$g(Q_3, Q_1)$
88	79	86	85	86

Table 4.1: *Percentage coverage of the true parameter values by approximate 95% credible intervals based on variational Bayes approximate posterior density functions. The percentages are based on 1000 replications.*

MCMC	MFVB
(533.68, 534.18)	(4.05, 4.52)

Table 4.2: *99% Wilcoxon confidence intervals based on run times in seconds for MCMC and MFVB fitting.*

MFVB is at least 100 times faster than our MCMC benchmark.

4.2.3 Application

We fitted (4.3) to geo-referenced data on sea-floor sediment pollution in the North Sea (source: Pebesma & Duin, 2005). The data are stored in the `pcb` data-frame within the R package `gstat` (Pebesma, 2004). The response variable is a measurement of polychlorinated biphenyl with Ballschmitter-Zell congener number 138 (PCB-138). The motivating study is concerned with spatial and temporal variability of PCB-138. For the purposes of illustration, we ignore the temporal aspect and focus on geographical variability in the mean and variance of the response. In the notation of model (4.2)–(4.3), the variables are $\mathbf{x} = (x_1, x_2)$ where

$x_1 = x$ -coordinate in the Universal Mobile Telecommunications System for Zone 31,

$x_2 = y$ -coordinate in the Universal Mobile Telecommunications System for Zone 31, and

$y =$ PCB-138 measured on the sediment fraction smaller than 63 parts per million, in $\mu\text{g}/\text{kg}$ dry matter.

The sample size is $n = 216$ and $K_u = K_v = 50$ thin plate spline bases were used for each functional fit. The estimated mean and standard deviation functions are shown in Figure 4.4. Both functions are seen to exhibit pronounced spatial effects. Simple geostatistical

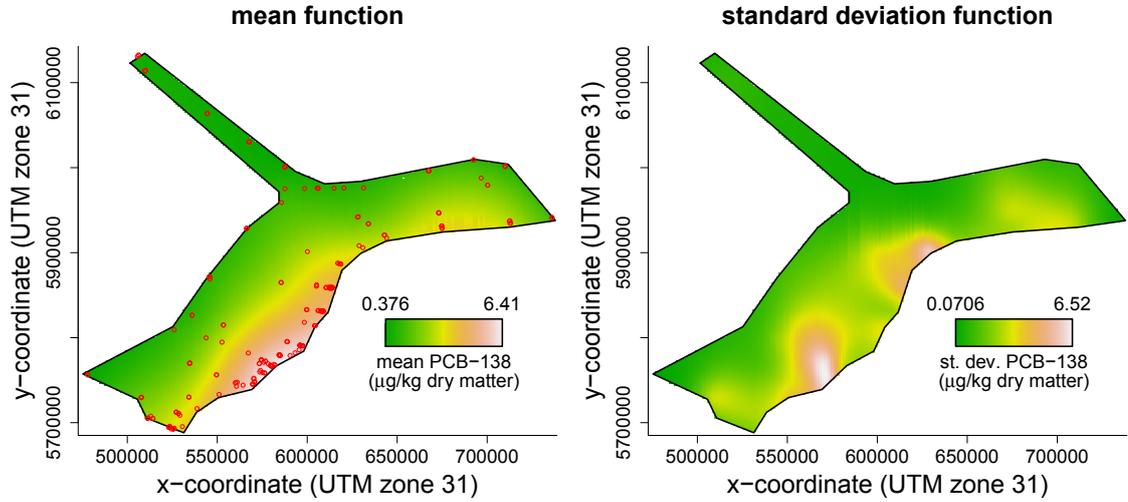


Figure 4.4: *Left panel: Fitted mean function for the polychlorinated biphenyl data described in the text, based on MFVB via Algorithm 2. Right panel: similar to top panel but for the standard deviation function.*

models that ignore the heteroscedasticity described by the right panel of Figure 4.4 would have erroneous prediction intervals.

Figure 4.5 gives a visual assessment of the convergence of the bivariate analogue of Algorithm 2 as it monitors successive values of $\log p(\mathbf{y}; q, \boldsymbol{\mu}_q(\boldsymbol{\omega}), \boldsymbol{\Sigma}_q(\boldsymbol{\omega}))$. It is evident that after approximately 75 iterations the Algorithm has reached convergence, after which inference was made to achieve Figure 4.4.

Higher dimensional heteroscedastic nonparametric regression can be achieved via Algorithm 2 with little notational change from the bivariate case treated here. The only required modification involves higher-dimensional thin plate spline basis functions instead of those given by (4.4).

4.3 Extension to additive models

The previous section showed how to construct tractable regression models for a bivariate set of continuous predictors modelled as a smooth function. However in many regression settings involving several continuous predictors, incorporation of multiple smooth functions may be more suited. The assumption made here is that of additivity, hence they are referred to as *additive models* (Ezekiel, 1924; Friedman & Stuetzle, 1981). Furthermore, we are dealing with models that need to accommodate heteroscedasticity, thus we

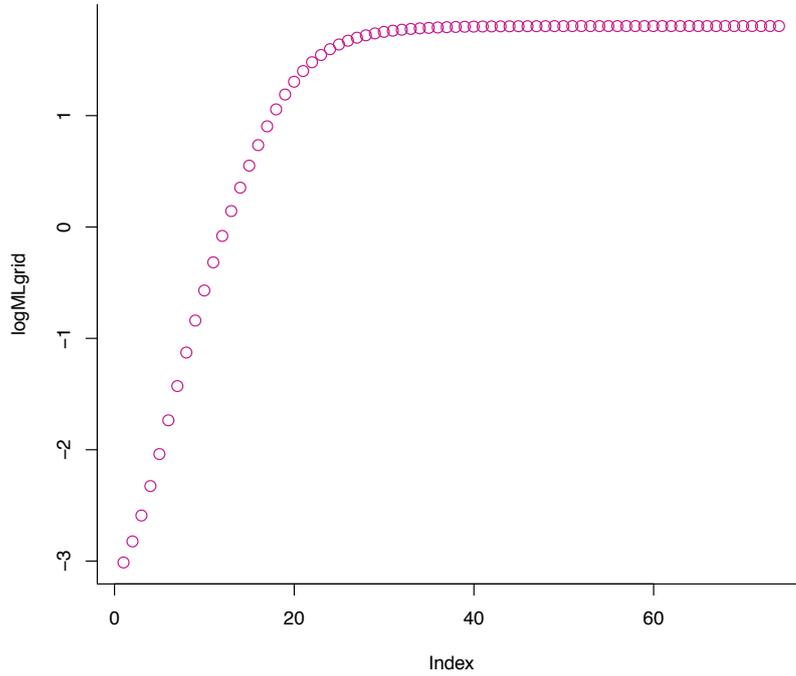


Figure 4.5: Successive values of $\log p(\mathbf{y}; q, \boldsymbol{\mu}_q(\boldsymbol{\omega}), \boldsymbol{\Sigma}_q(\boldsymbol{\omega}))$ for monitoring the convergence of non-conjugate MFVB given by the bivariate analogue of Algorithm 2.

incorporate the variance as a multiplicative function. Models of this type have a small literature with Rigby & Stasinopoulos (2005) being a key reference. Their generic form is

$$y_i \sim N\left(\beta_0 + \sum_{j=1}^d f_j(x_{ji}), \exp\left(\gamma_0 + \sum_{j=1}^d h_j(x_{ji})\right)\right), \quad 1 \leq i \leq n. \quad (4.6)$$

Here f_j and h_j , $1 \leq j \leq d$, are smooth but otherwise arbitrary functions. Penalised splines are easily extended to handle (4.6). Here we use penalized spline models of the form

$$\begin{aligned} f_j(x) &= \beta_j x + \sum_{k=1}^{K_j^u} u_{jk} z_{jk}^u(x), \quad u_{jk} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{u_j}^2) \\ \text{and } h_j(x) &= \gamma_j x + \sum_{k=1}^{K_j^v} v_{jk} z_{jk}^v(x), \quad v_{jk} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{v_j}^2) \end{aligned} \quad (4.7)$$

where the $\{z_{jk}^u : 1 \leq k \leq K_j^u\}$, $1 \leq j \leq d$, are spline basis functions of sizes K_j^u , analogue to those presented in Section 3.2. In addition the $\{z_{jk}^v : 1 \leq k \leq K_j^v\}$, $1 \leq j \leq d$, are similarly defined. The priors on the regression coefficients and standard deviation parameters are

$$\begin{aligned} \beta_j &\stackrel{\text{ind.}}{\sim} N(0, \sigma_\beta^2), \quad \gamma_j \stackrel{\text{ind.}}{\sim} N(0, \sigma_\gamma^2), \\ \sigma_{u_j} &\stackrel{\text{ind.}}{\sim} \text{Half-Cauchy}(A_u), \quad \sigma_{v_j} \stackrel{\text{ind.}}{\sim} \text{Half-Cauchy}(A_v), \quad 1 \leq j \leq d. \end{aligned} \quad (4.8)$$

The complete Bayesian hierarchical model corresponding to (4.6) is:

$$\begin{aligned}
 \mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{u}, \mathbf{v} &\sim \text{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_u\mathbf{u}, \text{diag}\{\exp(\mathbf{X}\boldsymbol{\gamma} + \mathbf{Z}_v\mathbf{v})\}), \\
 \mathbf{u}|\sigma_u^2 &\sim \text{N}\left(\mathbf{0}, \text{blockdiag}\left(\sigma_{\mathbf{u}_j}^2 \mathbf{I}_{K_{\mathbf{u}_j}}\right)_{1 \leq j \leq d}\right), \quad \mathbf{v}|\sigma_v^2 \sim \text{N}\left(\mathbf{0}, \text{blockdiag}\left(\sigma_{\mathbf{v}_j}^2 \mathbf{I}_{K_{\mathbf{v}_j}}\right)_{1 \leq j \leq d}\right), \\
 \sigma_{\mathbf{u}_j}^2|a_{\mathbf{u}_j} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\mathbf{u}_j}\right), \quad a_{\mathbf{u}_j} \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\mathbf{u}_j}^2\right), \\
 \sigma_{\mathbf{v}_j}^2|a_{\mathbf{v}_j} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\mathbf{v}_j}\right), \quad a_{\mathbf{v}_j} \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\mathbf{v}_j}^2\right), \\
 1 \leq j \leq d, \quad \boldsymbol{\beta} &\sim \text{N}(0, \sigma_\beta^2 \mathbf{I}_{d+1}), \quad \boldsymbol{\gamma} \sim \text{N}(0, \sigma_\gamma^2 \mathbf{I}_{d+1}).
 \end{aligned} \tag{4.9}$$

The differences of this model and that of Section 3.2 are that the design matrices are now

$$\begin{aligned}
 \mathbf{X} &\equiv [1 \ x_{1i} \ x_{2i} \ \cdots \ x_{di}]_{1 \leq i \leq n}, \\
 \mathbf{Z}_u &\equiv \left[\begin{array}{ccc} \left[z_k^{\mathbf{u}_1}(x_{1i}) \right]_{1 \leq i \leq n} & \left[z_k^{\mathbf{u}_2}(x_{2i}) \right]_{1 \leq i \leq n} & \cdots & \left[z_k^{\mathbf{u}_d}(x_{di}) \right]_{1 \leq i \leq n} \end{array} \right], \\
 \mathbf{Z}_v &\equiv \left[\begin{array}{ccc} \left[z_k^{\mathbf{v}_1}(x_{1i}) \right]_{1 \leq i \leq n} & \left[z_k^{\mathbf{v}_2}(x_{2i}) \right]_{1 \leq i \leq n} & \cdots & \left[z_k^{\mathbf{v}_d}(x_{di}) \right]_{1 \leq i \leq n} \end{array} \right]
 \end{aligned} \tag{4.10}$$

The coefficient vectors are

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, \quad \boldsymbol{\gamma} = \begin{bmatrix} \gamma_0 \\ \gamma_1 \\ \vdots \\ \gamma_d \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_d \end{bmatrix}, \quad \text{and} \quad \mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_d \end{bmatrix}, \tag{4.11}$$

where \mathbf{u}_j is the $K_{\mathbf{u}_j} \times 1$ vector containing the u_{jk} . It is again convenient to combine the mean function coefficients and variance function coefficients into single vectors as shown in (3.6). The Bayesian model given by (4.9) admits a closed form non-conjugate MFVB algorithm, with the regression coefficients for the full mean and variance functions each being Multivariate Normal. Algorithm 3 gives the details of these closed form updates.

Convergence of Algorithm 3 is assessed using the lower-bound of the marginal log-likelihood:

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_q(\boldsymbol{\omega}), \boldsymbol{\Sigma}_q(\boldsymbol{\omega})) &= \frac{1}{2} \left\{ \sum_{\ell=1}^d (K_{u_\ell} + K_{v_\ell}) + 4 \right\} - \frac{n}{2} \log(2\pi) + \log \Gamma\left(\frac{1}{2}(K_u + 1)\right) \\
 &\quad - 2 \log(\pi) + \log \Gamma\left(\frac{1}{2}(\sum_{\ell=1}^d K_{v_\ell} + d - 1)\right) - \sum_{\ell=1}^d \log(A_{u_\ell}) \\
 &\quad - \sum_{\ell=1}^d \log(A_{v_\ell}) - \frac{1}{2} \mathbf{1}^T (\mathbf{C}_\omega \boldsymbol{\mu}_q(\boldsymbol{\omega})) - \frac{1}{2} \mathbf{1}^T \left\{ \boldsymbol{\mu}_q(r_{v_\ell}^2) \odot \exp(\boldsymbol{\psi}_q(\boldsymbol{\omega})) \right\} \\
 &\quad - \frac{1}{2}(d+1) \log(\sigma_\beta^2) - \frac{1}{2}(d+1) \log(\sigma_\gamma^2) + \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\nu})| \\
 &\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\omega})| - \frac{1}{2\sigma_\beta^2} \left(\|\boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\beta})) \right) \\
 &\quad - \frac{1}{2\sigma_\gamma^2} \left\{ \|\boldsymbol{\mu}_q(\boldsymbol{\gamma})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\gamma})) \right\} - \frac{1}{2} (\sum_{\ell=1}^d K_{u_\ell} + 1) \log(B_q(\sigma_{u_\ell}^2)) \\
 &\quad - \frac{1}{2} (\sum_{\ell=1}^d K_{v_\ell} + 1) \log(B_q(\sigma_{v_\ell}^2)) - \log\left(\mu_q(1/\sigma_{u_\ell}^2) + A_{u_\ell}^{-2}\right) \\
 &\quad - \log\left(\mu_q(1/\sigma_{v_\ell}^2) + A_{v_\ell}^{-2}\right) + \mu_q(1/\sigma_{u_\ell}^2) \mu_q(1/a_{u_\ell}) + \mu_q(1/\sigma_{v_\ell}^2) \mu_q(1/a_{v_\ell}).
 \end{aligned}$$

Set up initial values:

$\boldsymbol{\mu}_{q(\omega)}$ a $(\sum_j^d K_{v_j} + d + 1) \times 1$ vector, $\boldsymbol{\Sigma}_{q(\omega)}$ a $(\sum_j^d K_{v_j} + d + 1) \times (\sum_j^d K_{v_j} + d + 1)$ positive definite matrix, $\mu_{q(1/\sigma_{u_j}^2)}, \mu_{q(1/\sigma_{v_j}^2)} > 0$, and $\boldsymbol{\mu}_{q(\mathbf{r}_\nu^2)}$ an $n \times 1$ vector.

Cycle through:

$$\begin{aligned} \boldsymbol{\psi}_{q(\omega)} &\leftarrow \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^\top \right) \right\} \\ \boldsymbol{\Sigma}_{q(\nu)} &\leftarrow \left(\mathbf{C}_\nu^\top \text{diag} \{ \boldsymbol{\psi}_{q(\omega)} \} \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_3 & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left(\mu_{q(1/\sigma_{u_j}^2)} I_{K_j}^u \right) \end{bmatrix} \right)^{-1} \\ \boldsymbol{\mu}_{q(\nu)} &\leftarrow \boldsymbol{\Sigma}_{q(\nu)} \mathbf{C}_\nu^\top \text{diag} \{ \boldsymbol{\psi}_{q(\omega)} \} \mathbf{y} \\ \boldsymbol{\Sigma}_{q(\omega)} &\leftarrow \left(\mathbf{C}_\omega^\top \text{diag} \{ \boldsymbol{\mu}_{q(\mathbf{r}_\nu^2)} \odot \boldsymbol{\psi}_{q(\omega)} \} \mathbf{C}_\omega + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_3 & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left(\mu_{q(1/\sigma_{v_j}^2)} I_{K_j}^v \right) \end{bmatrix} \right)^{-1} \\ \boldsymbol{\mu}_{q(\omega)} &\leftarrow \boldsymbol{\mu}_{q(\omega)} + \boldsymbol{\Sigma}_{q(\omega)} \left\{ \mathbf{C}_\omega^\top \left(\boldsymbol{\mu}_{q(\mathbf{r}_\nu^2)} \odot \boldsymbol{\psi}_{q(\omega)} - \mathbf{1} \right) \right. \\ &\quad \left. - \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_3 & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left(\mu_{q(1/\sigma_{v_j}^2)} I_{K_j}^v \right) \end{bmatrix} \boldsymbol{\mu}_{q(\omega)} \right\} \end{aligned}$$

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

$$\begin{aligned} \mu_{q(1/a_{u_j})} &\leftarrow 1 / \left(\mu_{q(1/\sigma_{u_j}^2)} + A_{u_j}^2 \right) ; \quad \mu_{q(1/a_{v_j})} \leftarrow 1 / \left(\mu_{q(1/\sigma_{v_j}^2)} + A_{v_j}^2 \right) \\ \mu_{q(1/\sigma_{u_j}^2)} &\leftarrow \frac{K_{u_j} + 1}{2\mu_{q(1/a_{u_j})} + \|\boldsymbol{\mu}_{q(u_j)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(u_j)})} \\ \mu_{q(1/\sigma_{v_j}^2)} &\leftarrow \frac{K_{v_j} + 1}{2\mu_{q(1/a_{v_j})} + \|\boldsymbol{\mu}_{q(v_j)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(v_j)})} \end{aligned}$$

until the absolute relative change in $\log p(\mathbf{y}; q, \boldsymbol{\mu}_{q(\omega)}, \boldsymbol{\Sigma}_{q(\omega)})$ is negligible.

Algorithm 3: *MFVB algorithm for the determination of the optimal parameters in $q^*(\omega)$, $q^*(\nu)$, $q^*(\sigma_{u_j}^2)$, $q^*(\sigma_{v_j}^2)$, $q^*(a_{u_j})$ and $q^*(a_{v_j})$, $1 \leq j \leq d$.*

The derivations of the optimal parameters in Algorithm 3 and the calculations corresponding to the lower-bound of the marginal log-likelihood are similar to that described in the Appendix 3.A and 3.B.

4.3.1 Application

Algorithm 3 has been applied to a data set from the Californian air pollution study described in Breiman & Friedman (1985). The data set consists of 345 observations on the four variables $x_{i1}, x_{i2}, x_{i3}, y_i$, $1 \leq i \leq 345$. The predictors are defined as

x_1 = pressure gradient (mm Hg) from Los Angeles International Airport to Daggett, California,

x_2 = inversion base height (feet)

x_3 = inversion base temperature (degrees Fahrenheit),

and the response variable is

y = ozone concentration (ppm) at Sandburg Air Force Base.

Our model is

$$y_i \stackrel{\text{ind.}}{\sim} N(\beta_0 + f_1(x_{1i}) + f_2(x_{2i}) + f_3(x_{3i}), \exp\{\gamma_0 + h_1(x_{1i}) + h_2(x_{2i}) + h_3(x_{3i})\}). \quad (4.12)$$

The data was standardised and the hyperparameters were set to have values $\sigma_\beta = \sigma_\gamma = A_u = A_v = 10^5$, to achieve non-informativity. After fitting, the resulting values were transformed to their original units. In addition the spline basis function sizes for modelling f_1 , f_2 and f_3 where 18 each. The same values were used for modelling h_1 , h_2 and h_3 . The approximation of the mean function and standard deviation function is shown in Figure 4.6. Inference based on MCMC was carried out using the same convergence criteria as explained in Section 4.2.2.

The estimated mean function f_1 has been vertically aligned to correspond to the response data by plotting, for \mathbf{x}_1 , the response with each of \mathbf{x}_2 and \mathbf{x}_3 set at their average. In other words, we have evaluated the estimate of f_2 at \bar{x}_2 and estimate of f_3 at \bar{x}_3 in order to vertically align the estimated f_1 curve. Similar alignments have been used for f_2 , f_3 , $\exp(h_1/2)$, $\exp(h_2/2)$, and $\exp(h_3/2)$.

As is seen from Figure 4.6 MFVB and MCMC are showing excellent agreement in

4.3. EXTENSION TO ADDITIVE MODELS

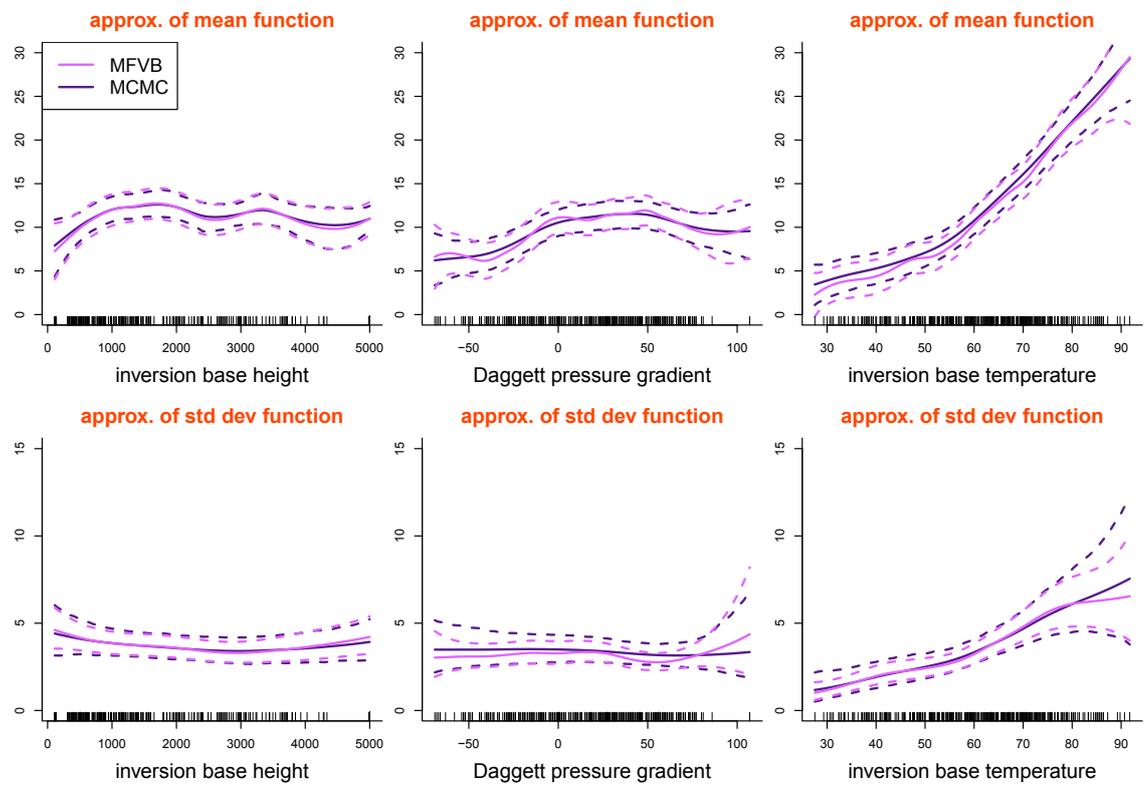


Figure 4.6: Top panels: approximation to the posterior mean functions f_1 , f_2 , and f_3 and their 95% credible sets. Bottom panels: approximation to the posterior standard deviation functions $\exp(h_1/2)$, $\exp(h_2/2)$, and $\exp(h_3/2)$ and their 95% credible sets.

their fits. The MFVB approach was markedly faster for this example. Even though the existence of heteroscedasticity in *inversion base height* and *Dagget pressure gradient* is somewhat low, there is noticeable heteroscedasticity in *inversion base temperature* which has been captured by model (4.12).

4.4 Real-time heteroscedastic nonparametric regression

The constant expansion of technological advancements in the 2010s so far has generated larger volumes of data than ever before. Not only is this information growing bigger but also increasing with higher speed.

Almost all nonparametric regression methodology presented to date make the assumption that the data are processed in batch, that is, at the same time. However, some disadvantages of batch processing include the requirement that analysis must wait until the entire data set has been collected. In the real-time case, the analysis updates once each new data point is collected. This is beneficial, and sometimes essential, for both high volume and/or velocity data. Many procedures used in single predictor nonparametric regression analysis use fully automated batch procedures however do not yet have an online/real-time modification. Recently, Luts *et al.* (2013) developed MFVB-based algorithms for performing semiparametric regression analysis in real time.

Algorithms that deal with the arrival of a stream of data points or just one data point per iteration have recently been developed for variational Bayesian inference by the machine learning community (e.g. Hoffman *et al.*, 2010; Wang *et al.*, 2011). This methodology has become known as *online mean field variational Bayes* or *online variational Bayes* for short. While most procedures to date involve storing a small amount of data in memory, they still need knowledge of the amount of data points from the beginning of the algorithm. In contrast, Luts *et al.* (2013) develop an online algorithm that uses past data only in the form of sufficient statistics and is relevant to fitting mixed models. We extend their approach by developing a variation of Algorithm 2 that allows for real-time heteroscedastic nonparametric regression. Unlike Luts *et al.* (2013) however, our real time MFVB algorithm is hindered by the fact that it requires storage of the entire data set. An extension to our algorithm, that only uses storage of the sufficient statistics would be worthwhile for future investigation.

The dependence on the data in Algorithm 2 is evident only in the response vector \mathbf{y}

4.4. REAL-TIME HETEROSCEDASTIC NONPARAMETRIC REGRESSION

and the design matrices \mathbf{C}_ν and \mathbf{C}_ω . These quantities are simple to update when a new observation y_{new} and its corresponding $\{(2 + K_u) \times 1\}$ and $\{(2 + K_v) \times 1\}$ vectors of design components $\mathbf{c}_{\nu,\text{new}}$ and $\mathbf{c}_{\omega,\text{new}}$ arrives. For instance, the new \mathbf{C}_ν matrix is

$$\mathbf{C}_\nu \leftarrow [\mathbf{C}_\nu^\top \ \mathbf{c}_{\nu,\text{new}}]^\top.$$

The corresponding updates for y_{new} and \mathbf{C}_ω are given in Algorithm 4. The starting values

1. Use Algorithm 2 to perform batch-based tuning runs, analogous to those described in Algorithm 2' of Luts *et al.* (2014), and determine a warm-up sample size n_{warm} for which convergence is validated.
2. Set $\boldsymbol{\mu}_{q(\nu)}$, $\boldsymbol{\Sigma}_{q(\nu)}$, $\boldsymbol{\mu}_{q(\omega)}$, $\boldsymbol{\Sigma}_{q(\omega)}$, $\boldsymbol{\mu}_{q(1/\sigma_u^2)}$, and $\boldsymbol{\mu}_{q(1/\sigma_v^2)}$ to their values obtained in the warm up batch-based tuning run with sample size n_{warm} . Next set y_{warm} to be the response vector on the first n_{warm} observations. Also set $\mathbf{C}_{\nu,\text{warm}}$ and $\mathbf{C}_{\omega,\text{warm}}$ to be the design matrices based on the first n_{warm} observations. Lastly assign $n \leftarrow n_{\text{warm}}$.
3. Cycle:

Read in y_{new} (1×1), $\mathbf{c}_{\nu,\text{new}}$ $\{(2 + K_u) \times 1\}$ and $\mathbf{c}_{\omega,\text{new}}$ $\{(2 + K_v) \times 1\}$; $n \leftarrow n + 1$

$$\mathbf{C}_\nu \leftarrow [\mathbf{C}_\nu^\top \ \mathbf{c}_{\nu,\text{new}}]^\top \quad ; \quad \mathbf{C}_\omega \leftarrow [\mathbf{C}_\omega^\top \ \mathbf{c}_{\omega,\text{new}}]^\top \quad ; \quad \mathbf{y} \leftarrow [\mathbf{y}^\top \ y_{\text{new}}]^\top$$

$$\boldsymbol{\mu}_{q(r_\nu^2)} \leftarrow \text{diagonal} \left\{ (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\mu}_{q(\nu)}) (\mathbf{y} - \mathbf{C}_\nu \boldsymbol{\mu}_{q(\nu)})^\top + \mathbf{C}_\nu \boldsymbol{\Sigma}_{q(\nu)} \mathbf{C}_\nu^\top \right\}$$

$$\psi_{q(\omega)} \leftarrow \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \text{diagonal} (\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^\top) \right\}$$

$$\boldsymbol{\Sigma}_{q(\nu)} \leftarrow \left(\mathbf{C}_\nu^\top \text{diag} (\psi_{q(\omega)}) \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_u^2)} \mathbf{I}_{K_u} \end{bmatrix} \right)^{-1}$$

$$\boldsymbol{\mu}_{q(\nu)} \leftarrow \boldsymbol{\Sigma}_{q(\nu)} \mathbf{C}_\nu^\top \text{diag} (\psi_{q(\omega)}) \mathbf{y}$$

$$\boldsymbol{\Sigma}_{q(\omega)} \leftarrow \left(\mathbf{C}_\omega^\top \text{diag} (\boldsymbol{\mu}_{q(r_\nu^2)} \odot \psi_{q(\omega)}) \mathbf{C}_\omega + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I}_{K_v} \end{bmatrix} \right)^{-1}$$

$$\boldsymbol{\mu}_{q(\omega)} \leftarrow \boldsymbol{\mu}_{q(\omega)} + \boldsymbol{\Sigma}_{q(\omega)} \left\{ \mathbf{C}_\omega^\top (\boldsymbol{\mu}_{q(r_\nu^2)} \odot \psi_{q(\omega)} - \mathbf{1}) - \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I}_{K_v} \end{bmatrix} \boldsymbol{\mu}_{q(\omega)} \right\}$$

$$\mu_{q(1/a_u)} \leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \quad ; \quad \mu_{q(1/a_v)} \leftarrow 1 / (\mu_{q(1/\sigma_v^2)} + A_v^{-2})$$

$$\mu_{q(1/\sigma_u^2)} \leftarrow (K_u + 1) / \{2\mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(u)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(u)})\}$$

$$\mu_{q(1/\sigma_v^2)} \leftarrow (K_v + 1) / \{2\mu_{q(1/a_v)} + \|\boldsymbol{\mu}_{q(v)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(v)})\}$$

until the analysis is complete or data are no longer available.

Algorithm 4: *MFVB algorithm for real-time determination of the optimal parameters in $q^*(\boldsymbol{\omega})$, $q^*(\boldsymbol{\nu})$, $q^*(\sigma_u^2)$, $q^*(\sigma_v^2)$, $q^*(a_u)$ and $q^*(a_v)$.*

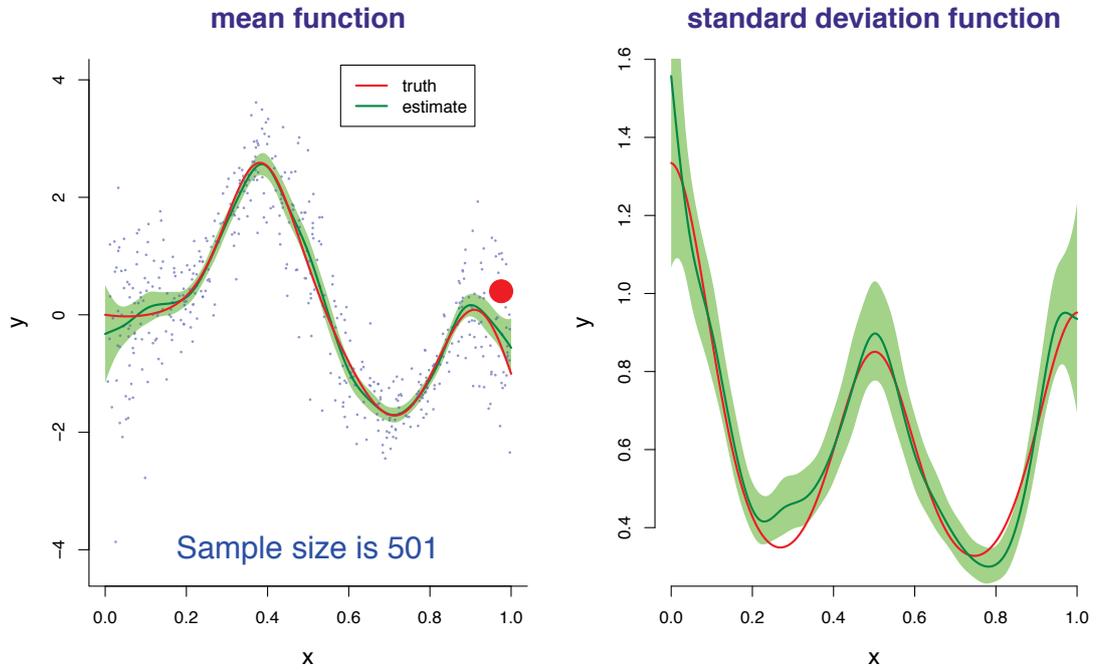
for the real-time procedure are determined by performing a sufficiently large batch fit.

When implementing Algorithm 4, the approximate posterior density functions of the model parameters are updated continually as each new data point arrives.

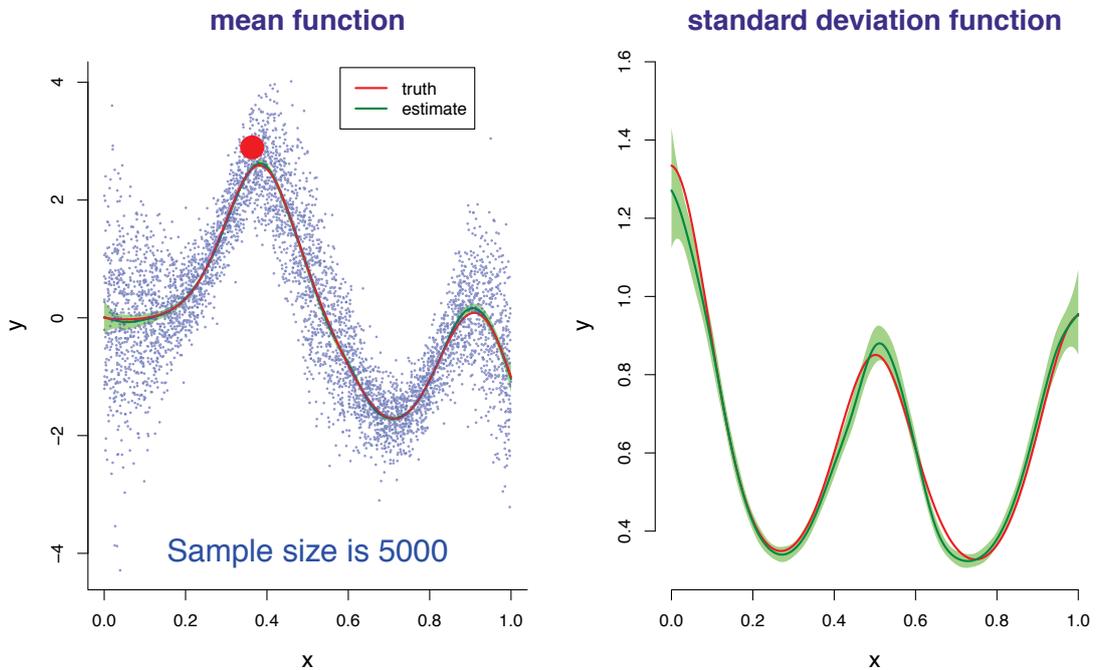
The web-site `realtime-semiparametric-regression.net` features a movie that illustrates Algorithm 4 for data simulated according to setting D from Table 3.1. The warm-up sample size is $n_{\text{warm}} = 500$. The link for the movie is titled `Heteroscedastic nonparametric regression` and portrays the effectiveness of real time processing for mean and standard deviation function fitting. Figure 4.7 illustrates the mean and variance function fits for the first online data point at $n = 501$ and again at the last online data point where $n = 5000$. As each new data point is processed, the mean and variance function fits improve as seen from this comparison.

4.5 Discussion

This chapter extends the work in Chapter 3 as we consider a wider class of models for fitting and inference using non-conjugate MFVB. We have developed closed form algorithms for both fast batch and real-time fitting and inference for a range of semiparametric regression models that have a heteroscedastic component. The methodology developed here applies to increasingly complex models as a result of the locality property of MFVB inference methods in general. In both the simulated and actual data settings, the new methodology is seen to perform very well.



(a) Resulting mean and standard deviation function fits after first new data point arrives.



(b) Resulting mean and standard deviation function fits once all 5000 data points have arrived.

Figure 4.7: Real-time approximate non-conjugate MFVB regression fits for the simulated data according to setting D from Table 3.1.

Chapter 5

Mean field variational Bayes for group-specific curve models

5.1 Introduction

In this chapter we develop MFVB algorithms catered to fitting large data sets that exhibit multilevel and longitudinal structures. In addition to developing the general MFVB algorithms that we have so far seen in the previous chapters of this thesis, here we also look into streamlining these algorithms in terms of storage and number of operations needed. Streamlining of such algorithms is essential when using MFVB for fitting and inference in multilevel and longitudinal models. This is often the case since a naïve implementation would result in cubic dependence on the number of groups within a level of the model being considered.

Lee & Wand (2015) develop streamlined algorithms for treatment of two-level models with Gaussian and binary responses and show that the number of operations required are linear in the number of groups in each level. This allows for much faster and highly accurate Bayesian inference for large longitudinal and multilevel models. We adopt the methodology used in Lee & Wand (2015) to streamline the two-level MFVB algorithm used in Section 5.2.3 and also extend their methodology to cater to the three-level model considered in Section 5.3.2. These algorithms have some inferential inaccuracy, however, as represented by the diagnostic illustrations in this chapter, these inaccuracies are relatively insignificant when compared to the computational time and storage savings afforded by the streamlined MFVB algorithm. The software package `Infer.NET` allows treatment of

such models, however, its non-streamlined implementation makes MFVB very slow when the number of groups is large.

Multilevel and longitudinal models, e.g. Diggle *et al.* (2002), Fitzmaurice *et al.* (2012), Gelman & Hill (2006) and Goldstein (2011), is a prominent area in Statistics, however the union of these areas and that of variational methods has not yet appeared to be eminent. Recent contributions are found in Ormerod & Wand (2010) and Luts *et al.* (2014), which use MFVB for longitudinal/multilevel data but invert the effects covariance matrix parameter updates in a naïve fashion. The main idea of this chapter is to produce MFVB algorithms for multilevel/longitudinal data based on the naïve inversion of the effects covariance matrix and compare this to the algorithms using the streamlined inversion of the effects covariance matrix. We restrict attention to the two-level and three-level Gaussian response models.

The motivation for using a streamlined MFVB algorithm in this chapter comes from a three-level longitudinal data set that represents frequency dependent backscatter coefficients for induced tumors in rodents (Simpson, 2013). The structure of this data led to the consideration of a variant of the general MFVB approach, referred to as streamlined MFVB, which we explore in further detail in this chapter.

Section 5.2 gives description of the two-level Gaussian response model in terms of its construction and corresponding MFVB methodology and Section 5.3 provides similar descriptions for the three-level Gaussian response model. Comprehensive simulation studies and real data examples are considered for both the two-level and three-level Gaussian response models. Appendix 5.A and 5.B provide the derivations for the naïve two-level model Algorithm and its corresponding lower bound on the marginal log-likelihood.

5.2 Two-level Gaussian response model

We consider the development of a model for fitting longitudinal and multilevel data such as that illustrated in Figure 5.1. There are a number of things to consider when developing a model that best represents data of this type. For instance, the repeated measures or longitudinal aspect is crucial since one would expect, upon inspection of Figure 5.1, that utilising *within-subject* information would be helpful. So, at the very least we would want to include a variance component that controls for the deviation within each subject, i.e., a random intercept. In addition, one might also consider a random slope which would

5.2. TWO-LEVEL GAUSSIAN RESPONSE MODEL

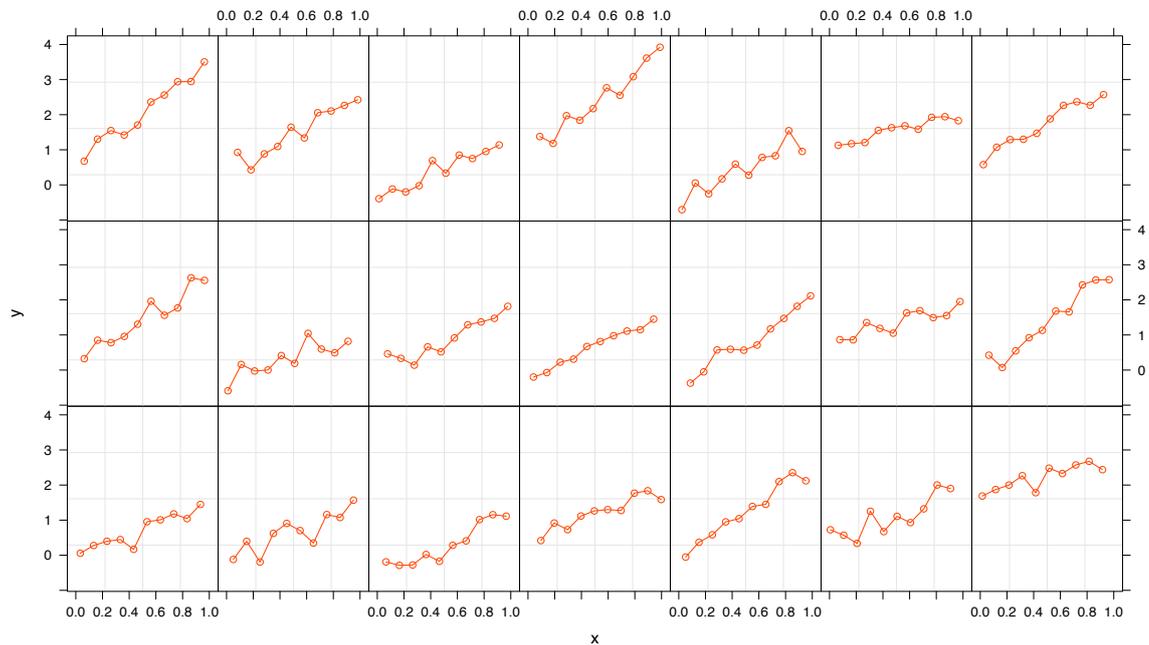


Figure 5.1: *Example of a longitudinal two-level model data-set.*

account for possible variability in the slopes of the growth curves.

The nonlinearity of such data suggests that incorporation of a nonparametric function f would be beneficial. However, f would represent a global function which does not depend on each subject. In order to cater for each individual's nonlinear component, we might consider a model of the following form:

$$y_{ij} = f(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad (5.1)$$

where f is a smooth function and the g_i , $1 \leq i \leq m$, are random functions with mean zero. Both f and g can be modeled as regression splines and then placed into the mixed model framework.

The regression splines associated with f have a fixed component that caters for the linear part of the model, and a random component that allows for departures from linearity. It is important to note, on the other hand, since the g_i are random with zero mean, the linear part of the regression spline is random, rather than fixed. A linear penalised spline

model for (5.1) is

$$f(x_{ij}) = \beta_0 + \beta_1 x_{ij} + \sum_{\ell=1}^{L^{\text{gbl}}} u_{\ell}^{\text{gbl}} z_{\ell}^{\text{gbl}}(x_{ij}) \quad (5.2)$$

$$\text{and } g_i(x_{ij}) = \delta_{0i} + \delta_{1i} x_{ij} + \sum_{\ell=1}^{L^{\text{grp}}} u_{i\ell}^{\text{grp}} z_{\ell}^{\text{grp}}(x_{ij}), \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i,$$

where the u_{ℓ}^{gbl} and $z_{\ell}^{\text{gbl}}(x_{ij})$ are the spline coefficient and spline basis function for the ℓ th knot in the global part of the model, that is, corresponding to the highest level of the two-level model structure. The $u_{i\ell}^{\text{grp}}$ and $z_{\ell}^{\text{grp}}(x_{ij})$ are the spline coefficient and spline basis function for the ℓ th knot of the i th subject, that corresponding to the lower level of the two-level model structure. A popular choice for the z_{ℓ}^{gbl} and z_{ℓ}^{grp} are suitably transformed cubic O'Sullivan splines, as described in Section 1.10. Equations (5.1) and (5.2) can be represented in the mixed model representation as follows.

5.2.1 Mixed model representation

We define the following linear and non-linear predictors:

$$\mathbf{x}_{ij} = [1 \ x_{ij}], \quad \mathbf{z}_{ij}^{\text{gbl}} = [z_1^{\text{gbl}}(x_{ij}) \ \dots \ z_{L^{\text{gbl}}}^{\text{gbl}}(x_{ij})] \quad \text{and} \quad \mathbf{z}_{ij}^{\text{grp}} = [z_1^{\text{grp}}(x_{ij}) \ \dots \ z_{L^{\text{grp}}}^{\text{grp}}(x_{ij})].$$

In addition, we define

$$\begin{aligned} \boldsymbol{\beta} &= [\beta_0 \ \beta_1]^{\top}, \quad \boldsymbol{\delta}_i = [\delta_{0i} \ \delta_{1i}]^{\top}, \\ \mathbf{u}^{\text{gbl}} &= [u_1^{\text{gbl}} \ \dots \ u_{L^{\text{gbl}}}^{\text{gbl}}]^{\top}, \quad \mathbf{u}_i^{\text{grp}} = [u_{i1}^{\text{grp}} \ \dots \ u_{iL^{\text{grp}}}^{\text{grp}}]^{\top}. \end{aligned} \quad (5.3)$$

We can now write (5.2) as

$$\begin{aligned} f(x_{ij}) &= \mathbf{x}_{ij} \boldsymbol{\beta} + \mathbf{z}_{ij}^{\text{gbl}} \mathbf{u}^{\text{gbl}}, \\ g_i(x_{ij}) &= \mathbf{x}_{ij} \boldsymbol{\delta}_i + \mathbf{z}_{ij}^{\text{grp}} \mathbf{u}_i^{\text{grp}}. \end{aligned}$$

We set $\text{Cov}(\mathbf{u}^{\text{gbl}}) = \sigma_{\text{gbl}}^2 \mathbf{I}$, $\text{Cov}(\boldsymbol{\delta}_i) = \boldsymbol{\Sigma}$, an unstructured 2×2 covariance matrix which allows for deviation from the typical linear component and lastly we set $\text{Cov}(\mathbf{u}_i^{\text{grp}}) = \sigma_{\text{grp}}^2 \mathbf{I}$.

We can now incorporate (5.1) into a mixed model framework by letting:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{11} \\ \vdots \\ \mathbf{x}_{1n_1} \\ \vdots \\ \mathbf{x}_{m1} \\ \vdots \\ \mathbf{x}_{mn_m} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \text{and} \quad (5.4)$$

$$\mathbf{Z} = \begin{bmatrix} z_{11}^{\text{gbl}} & \mathbf{x}_{11} & z_{11}^{\text{grp}} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ z_{1n_1}^{\text{gbl}} & \mathbf{x}_{1n_1} & z_{1n_1}^{\text{grp}} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ z_{21}^{\text{gbl}} & \mathbf{0} & \mathbf{0} & \mathbf{x}_{21} & z_{21}^{\text{grp}} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ z_{2n_2}^{\text{gbl}} & \mathbf{0} & \mathbf{0} & \mathbf{x}_{2n_2} & z_{2n_2}^{\text{grp}} & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ z_{m1}^{\text{gbl}} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{x}_{m1} & z_{m1}^{\text{grp}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ z_{mn_m}^{\text{gbl}} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{x}_{mn_m} & z_{mn_m}^{\text{grp}} \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}^{\text{gbl}} \\ \boldsymbol{\delta}_1 \\ \mathbf{u}_1^{\text{grp}} \\ \boldsymbol{\delta}_2 \\ \mathbf{u}_2^{\text{grp}} \\ \vdots \\ \boldsymbol{\delta}_m \\ \mathbf{u}_m^{\text{grp}} \end{bmatrix}, \quad (5.5)$$

so we can use the mixed model representation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon},$$

$$\text{where } \text{Cov} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\varepsilon} \end{bmatrix} = \begin{bmatrix} \sigma_{\text{gbl}}^2 \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\boldsymbol{\Sigma}, \sigma_{\text{grp}}^2 \mathbf{I}_{L_{\text{grp}}}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \sigma_{\boldsymbol{\varepsilon}}^2 \mathbf{I} \end{bmatrix}. \quad (5.6)$$

This semiparametric regression model can be treated using the following Gaussian linear mixed model:

$$\mathbf{y}|\mathbf{u} \sim \text{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_{\boldsymbol{\varepsilon}}^2 \mathbf{I}), \quad \mathbf{u} \sim \text{N}\left(\mathbf{0}, \begin{bmatrix} \sigma_{\text{gbl}}^2 \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\boldsymbol{\Sigma}, \sigma_{\text{grp}}^2 \mathbf{I}_{L_{\text{grp}}}) \end{bmatrix}\right). \quad (5.7)$$

5.2.2 Bayesian inference

Fitting (5.7) using standard mixed model software such as `lme()` (Pinheiro *et al.*, 2009) in the R computing language is achievable, but time consuming. Alternatively, we can work with a hierarchical Bayesian version of (5.7) which warrants direct implementation in standard software. We enforce σ_{gbl} , the $\sigma_{\text{grp},i}$, $1 \leq i \leq m$, and the square-rooted diagonal entries of $\boldsymbol{\Sigma}$ to possess Half-t prior distributions. This specification also allows for corre-

lation parameters within Σ to have uniform distributions on $(-1, 1)$, as shown in Huang *et al.* (2013). Such priors are attained through auxiliary variable constructions exhibited in Result 1.4.6.

In its entirety, the Bayesian hierarchical two-level Gaussian response model we deal with here is:

$$\begin{aligned}
 \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim \text{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), & \boldsymbol{\beta} &\sim \text{N}(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_2), \\
 \mathbf{u} | \sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2 &\sim \text{N}\left(\mathbf{0}, \begin{bmatrix} \sigma_{\text{gbl}}^2 \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\Sigma, \sigma_{\text{grp}}^2 \mathbf{I}_{L_{\text{grp}}}) \end{bmatrix}\right), \\
 \sigma_\varepsilon^2 | a_\varepsilon &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_\varepsilon\right), & a_\varepsilon &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_\varepsilon^2\right), \\
 \sigma_{\text{gbl}}^2 | a_{\text{gbl}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{gbl}}\right), & a_{\text{gbl}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{gbl}}^2\right), \\
 \Sigma | a_{\Sigma,1}, a_{\Sigma,2} &\sim \text{Inverse-Wishart}(\nu + 2 - 1, 2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})), \\
 a_{\Sigma,j} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\Sigma,j}^2\right), & 1 \leq j \leq 2, \\
 \sigma_{\text{grp}}^2 | a_{\text{grp}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{grp}}\right), & a_{\text{grp}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{grp}}^2\right).
 \end{aligned} \tag{5.8}$$

where a_ε , a_{gbl} , $a_{\Sigma,j}$, and a_{grp} represent the auxiliary variables that impose the aforementioned non-informative prior distributions on σ_ε^2 , σ_{gbl}^2 , Σ , and σ_{grp}^2 , respectively. We also note that (5.8) encompasses the setting of hyperparameters ν , σ_β , A_ε , A_{gbl} , $A_{\Sigma,1}$, $A_{\Sigma,2}$ and A_{grp} , all of which are bound to be positive. The recommended setting for ν is 2, which ensures that Property 4 of Huang *et al.* (2013) holds. Non-informativity is achieved for the remaining hyperparameters when they are set to very large values. Our default setting, assuming the data has been transformed to have mean zero and unit standard deviation, is 10^5 .

Figure 5.2 shows a directed acyclic graph representation of model (5.8) where $\mathbf{a}_\Sigma = [a_{\Sigma,1}, a_{\Sigma,2}]^\top$, $\boldsymbol{\delta} = [\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_m]^\top$ and $\mathbf{u}_{\text{grp}} = [\mathbf{u}_1^{\text{grp}}, \dots, \mathbf{u}_m^{\text{grp}}]^\top$.

5.2.3 Mean field variational Bayes methodology

The underpinning assumption of the MFVB approach for model (5.8) is to impose an approximate product density restriction to the joint posterior density function of the form

$$\begin{aligned}
 p(\boldsymbol{\beta}, \mathbf{u}, a_\varepsilon, a_{\text{gbl}}, \mathbf{a}_\Sigma, a_{\text{grp}}, \sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2 | \mathbf{y}) &\approx q(\boldsymbol{\beta}, \mathbf{u}, a_\varepsilon, a_{\text{gbl}}, \mathbf{a}_\Sigma, a_{\text{grp}}) \\
 &\times q(\sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2).
 \end{aligned} \tag{5.9}$$

Simplicity of the requisite calculations arise as a result of the graphical layout of model (5.8) as conveyed by Figure 5.2. The graphical structure of (5.8) is useful in determining

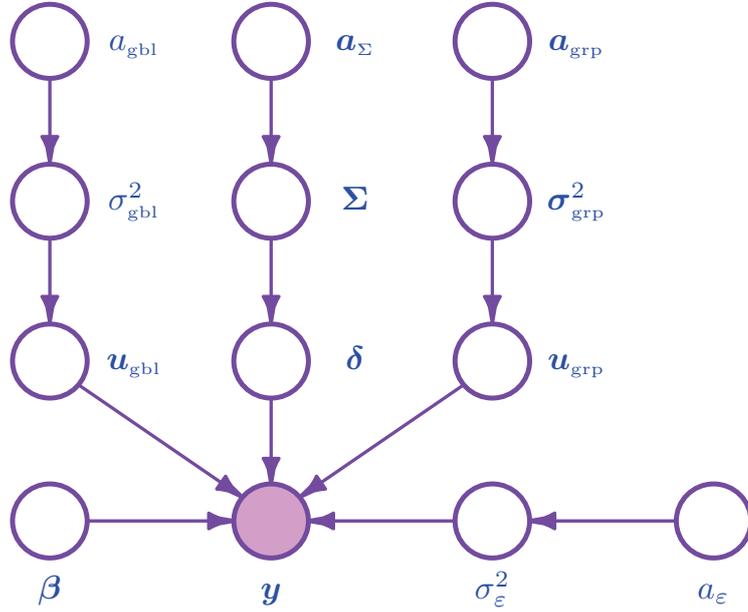


Figure 5.2: Directed acyclic graph corresponding to the Bayesian hierarchical two-level Gaussian response model in (5.8).

additional factorizations, i.e., *induced factorizations*. Factorizations such as these, can be observed using a straightforward graphical test, known as *moralisation*, as explained in Section 1.8.2. For example, all paths between σ_ε^2 and $\{\sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2\}$ in Figure 5.2 must visit at least one of $\{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\}$. That is, the set of nodes $\{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\}$ separates σ_ε^2 from the set $\{\sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2\}$. The use of Theorem 1.8.1 gives the result

$$\sigma_\varepsilon^2 \perp\!\!\!\perp \{\sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2\} \mid \{\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}\}.$$

Continuing this process for all nodes in Figure (5.2), we achieve the following induced factorization:

$$q(\boldsymbol{\beta}, \mathbf{u}, a_\varepsilon, a_{\text{gbl}}, a_\Sigma, a_{\text{grp}}, \sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \sigma_{\text{grp}}^2, \Sigma) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma_\varepsilon^2) q(\sigma_{\text{gbl}}^2) q(\Sigma) q(\sigma_{\text{grp}}^2) q(a_\varepsilon) q(a_{\text{gbl}}) \left\{ \prod_{j=1}^2 q(a_{\Sigma, j}) \right\} q(a_{\text{grp}}). \quad (5.10)$$

The q -densities are chosen to minimize the Kullback-Leibler distance between the left-hand-side and right-hand-side of (5.9). As outlined in Section 1.5 the optimal q -densities, denoted by q^* , can be obtained through an iterative scheme arising from relationships such as

$$q^*(\Sigma) \propto \exp E_q \{ \log p(\Sigma | \text{rest}) \}, \quad (5.11)$$

where ‘rest’ denotes the random variables in the model not including Σ . The full conditional distributions each have standard forms resulting in closed form expressions for the optimal q -densities. As a result, the optimal factors on the right-hand side of (5.10) have the following forms (details of the derivation can be found in Appendix 5.A):

$$\begin{aligned}
 q^*(\boldsymbol{\beta}, \mathbf{u}) & \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \text{ density function,} \\
 q^*(\sigma_\varepsilon^2) & \text{ is the Inverse-Gamma}\left(\frac{1}{2}\left(\sum_{i=1}^m n_i + 1\right), B_{q(\sigma_\varepsilon^2)}\right) \text{ density function,} \\
 q^*(\sigma_{\text{gbl}}^2) & \text{ is the Inverse-Gamma}\left(\frac{1}{2}(L_{\text{gbl}} + 1), B_{q(\sigma_{\text{gbl}}^2)}\right) \text{ density function,} \\
 q^*(\boldsymbol{\Sigma}) & \text{ is the Inverse-Wishart}(\nu + m + 1, B_{q(\boldsymbol{\Sigma})}) \text{ density function,} \\
 q^*(\sigma_{\text{grp}}^2) & \text{ is the Inverse-Gamma}\left(\frac{1}{2}(m L_{\text{grp}} + 1), B_{q(\sigma_{\text{grp}}^2)}\right) \text{ density function,} \\
 q^*(a_\varepsilon) & \text{ is the Inverse-Gamma}(1, B_{q(a_\varepsilon)}) \text{ density function,} \\
 q^*(a_{\text{gbl}}) & \text{ is the Inverse-Gamma}(1, B_{q(a_{\text{gbl}})}) \text{ density function,} \\
 q^*(a_{\Sigma, j}) & \text{ is the Inverse-Gamma}\left(\frac{\nu}{2} + 1, B_{q(a_{\Sigma, j})}\right) \text{ density function, } 1 \leq j \leq 2, \\
 q^*(a_{\text{grp}}) & \text{ is the Inverse-Gamma}(1, B_{q(a_{\text{grp}})}) \text{ density function,}
 \end{aligned} \tag{5.12}$$

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 rate parameter of $q^*(\sigma_\varepsilon^2)$, $B_{q(\sigma_{\text{gbl}}^2)}$, the rate parameter of $q^*(\sigma_{\text{gbl}}^2)$, $B_{q(\boldsymbol{\Sigma})}$, the rate matrix of $q^*(\boldsymbol{\Sigma})$, $B_{q(\sigma_{\text{grp}}^2)}$, the rate parameter of $q^*(\sigma_{\text{grp}}^2)$, $B_{q(a_\varepsilon)}$, the rate parameter of $q^*(a_\varepsilon)$, $B_{q(a_{\text{gbl}})}$, the rate parameter of $q^*(a_{\text{gbl}})$, $B_{q(a_{\Sigma, j})}$, the rate parameter of $q^*(a_{\Sigma, j})$ and $B_{q(a_{\text{grp}})}$, the rate parameter of $q^*(a_{\text{grp}})$. The parameters in these optimal q -densities are obtained through a naïve implementation which results in Algorithm 5 and uses the notation $\mathbf{C} \equiv [\mathbf{X} | \mathbf{Z}]$. As seen from (5.5) the \mathbf{Z} matrix requires storage of zero’s which can become computationally expensive given a large value for m . We aim to alleviate the restriction of storing these zero’s in Section 5.2.4.

Convergence of Algorithm 5 is determined using the marginal log-likelihood lower

Set up initial values:

$$\mu_{q(1/\sigma_\varepsilon^2)}, \mu_{q(1/\sigma_{\text{gbl}}^2)}, \mu_{q(1/\sigma_{\text{grp}}^2)}, \mu_{q(1/a_\varepsilon)}, \mu_{q(1/a_{\text{gbl}})}, \mu_{q(1/a_{\text{grp}})} > 0, \mu_{q(1/a_{\Sigma,j})} > 0,$$

$1 \leq j \leq 2$, $\Sigma_{q(\beta, \mathbf{u})}$, and $\mathbf{M}_{q(\Sigma^{-1})}$ positive definite.

Cycle through:

$$\begin{aligned} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} &\leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top \mathbf{y} \\ \Sigma_{q(\beta, \mathbf{u})} &\leftarrow \left(\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} \right. \\ &\quad \left. + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_{\text{gbl}}^2)} \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag} \left(\mathbf{M}_{q(\Sigma^{-1})}, \mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}} \right) \right] \right)^{-1} \end{aligned}$$

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

$$\mu_{q(1/\sigma_\varepsilon^2)} \leftarrow \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) / B_{q(\sigma_\varepsilon^2)}$$

$$B_{q(a_\varepsilon)} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1/B_{q(a_\varepsilon)}$$

$$B_{q(\sigma_{\text{gbl}}^2)} \leftarrow \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{gbl}})}) \right\} + \mu_{q(1/a_{\text{gbl}})}$$

$$\mu_{q(1/\sigma_{\text{gbl}}^2)} \leftarrow \frac{1}{2} (L_{\text{gbl}} + 1) / B_{q(\sigma_{\text{gbl}}^2)}$$

$$B_{q(a_{\text{gbl}})} \leftarrow \mu_{q(1/\sigma_{\text{gbl}}^2)} + A_{\text{gbl}}^{-2}; \quad \mu_{q(1/a_{\text{gbl}})} \leftarrow 1/B_{q(a_{\text{gbl}})}$$

$$B_{q(\sigma_{\text{grp}}^2)} \leftarrow \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grp}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{grp}})}) \right\} + \mu_{q(1/a_{\text{grp}})}$$

$$\mu_{q(1/\sigma_{\text{grp}}^2)} \leftarrow \frac{1}{2} (m L_{\text{grp}} + 1) / B_{q(\sigma_{\text{grp}}^2)}$$

$$B_{q(a_{\text{grp}})} \leftarrow \mu_{q(1/\sigma_{\text{grp}}^2)} + A_{\text{grp}}^{-2}; \quad \mu_{q(1/a_{\text{grp}})} \leftarrow 1/B_{q(a_{\text{grp}})}$$

$$\mathbf{B}_{q(\Sigma)} \leftarrow \sum_{i=1}^m \left\{ \boldsymbol{\mu}_{q(\delta_i)} \boldsymbol{\mu}_{q(\delta_i)}^\top + \Sigma_{q(\delta_i)} \right\} + 2\nu \text{diag} \left(\mu_{q(1/a_{\Sigma,1})}, \mu_{q(1/a_{\Sigma,2})} \right)$$

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

For $j = 1, 2$:

$$B_{q(a_{\Sigma,j})} \leftarrow \nu \left(\mathbf{M}_{q(\Sigma^{-1})} \right)_{jj} + 1/A_{\Sigma,j}^2, \quad \mu_{q(1/a_{\Sigma,j})} \leftarrow \left(\frac{\nu}{2} + 1 \right) / B_{q(a_{\Sigma,j})}$$

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

Algorithm 5: *Naïve MFVB algorithm for the estimation of the optimal parameters in $q^*(\beta, \mathbf{u})$, $q^*(\sigma_\varepsilon^2)$, $q^*(a_\varepsilon)$, $q^*(\sigma_{\text{gbl}}^2)$, $q^*(a_{\text{gbl}})$, $q^*(\sigma_{\text{grp}}^2)$, $q^*(a_{\text{grp}})$, $q^*(\Sigma)$, and $q^*(\mathbf{a}_\Sigma)$.*

bound (Appendix 5.B shows details on this derivation):

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2}(\nu + p - 1) \log(2\nu) - \frac{1}{2} \sum_{i=1}^m n_i \log(2\pi) - 4 \log(\pi) - \frac{p}{2} \log(\sigma_\beta^2) \\
 &\quad - \frac{1}{\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| \\
 &\quad + \frac{1}{2} \{p + L_{\text{gbl}} + m(p + L_{\text{grp}})\} - \log(C_{p, \nu+p-1}) + \log(C_{p, \nu+m+p-1}) \\
 &\quad - \frac{1}{2} (\nu + p + m - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma})}| + \log \Gamma \left\{ \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \right\} \\
 &\quad - \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \log(B_{q(\sigma_\varepsilon^2)}) + \log \Gamma \left\{ \frac{1}{2} (L_{\text{gbl}} + 1) \right\} \\
 &\quad - \frac{1}{2} (L_{\text{gbl}} + 1) \log \left(B_{q(\sigma_{\text{gbl}}^2)} \right) + \log \Gamma \left\{ \frac{1}{2} (m \times L_{\text{grp}} + 1) \right\} \\
 &\quad - \frac{1}{2} (m \times L_{\text{grp}} + 1) \log \left(B_{q(\sigma_{\text{grp}}^2)} \right) - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) \\
 &\quad + \mu_{q(1/\sigma_\varepsilon^2)} \mu_{q(1/a_\varepsilon)} - \log(A_{\text{gbl}}) - \log(B_{q(a_{\text{gbl}})}) + \mu_{q(1/\sigma_{\text{gbl}}^2)} \mu_{q(1/a_{\text{gbl}})} \\
 &\quad - \log(A_{\text{grp}}) - \log(B_{q(a_{\text{grp}})}) + \mu_{q(1/\sigma_{\text{grp}}^2)} \mu_{q(1/a_{\text{grp}})} - \sum_{j=1}^p \log(A_{\Sigma, j}) \\
 &\quad + p \log \Gamma \left\{ \frac{1}{2} (\nu + p) \right\} - \frac{1}{2} (\nu + p - 1) \sum_{j=1}^p \log(B_{q(a_{\Sigma, j})}) \\
 &\quad + \sum_{j=1}^p \nu \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma, j})},
 \end{aligned} \tag{5.13}$$

where $p = 2$ is the number of columns in the \mathbf{X} matrix, which also corresponds to the dimension of $\boldsymbol{\Sigma}$.

5.2.4 Streamlining Mean field variational Bayes for the two-level Gaussian response model

The naïve nature of the implementation given by Algorithm 5, stems from the fact that the update for the covariance matrix $\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}$ requires storing and inverting a matrix which has its row length equal to that of the \mathbf{C} matrix and thus can become infeasible for large multilevel/longitudinal datasets. Lee & Wand (2015) develop a streamlined MFVB algorithm for efficient fitting and inference in large longitudinal and multilevel data sets. In particular they develop streamlined MFVB algorithms for the two-level Gaussian response model and the two-level binary response model. We focus on the former algorithm and slightly adjust their Algorithm 2 to cater to (5.8).

In keeping with the notation used in Lee & Wand (2015) we further partition the

following vectors and matrices:

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

where

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \mathbf{X}_i = \begin{bmatrix} 1 & x_{i1} \\ \vdots & \vdots \\ 1 & x_{in_i} \end{bmatrix}, \quad \mathbf{Z}_i^{\text{gbl}} = \begin{bmatrix} z_{i1}^{\text{gbl}\top} \\ \vdots \\ z_{in_i}^{\text{gbl}\top} \end{bmatrix}, \quad \text{and } \mathbf{Z}_i^{\text{R}} = \begin{bmatrix} x_{i1}^\top & z_{i1}^{\text{grp}\top} \\ \vdots & \vdots \\ x_{in_i}^\top & z_{in_i}^{\text{grp}\top} \end{bmatrix}.$$

In addition, we partition the following global and random effects:

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}^{\text{gbl}} \\ \mathbf{u}^{\text{R}} \end{bmatrix}, \quad \text{and } \mathbf{Z} = [\mathbf{Z}^{\text{gbl}} \quad \mathbf{Z}^{\text{R}}],$$

where $\mathbf{u}^{\text{R}} = [(\mathbf{u}_1^{\text{R}})^\top \dots (\mathbf{u}_m^{\text{R}})^\top]^\top$, $\mathbf{u}_i^{\text{R}} = [\boldsymbol{\delta}_i^\top (\mathbf{u}_i^{\text{grp}})^\top]^\top$ and \mathbf{u}^{gbl} is defined in (5.3). Finally, Lee & Wand (2015) suggest the useful notation $\mathbf{C}^{\text{gbl}} = [\mathbf{X} \quad \mathbf{Z}^{\text{gbl}}]$ where $\mathbf{C}_i^{\text{gbl}} = [\mathbf{X}_i \quad \mathbf{Z}_i^{\text{gbl}}]$ is the sub-matrix of \mathbf{C}^{gbl} corresponding to the i th subject. With the aid of this new notation, a streamlined version of Algorithm 5 is presented in Algorithm 6.

Convergence of Algorithm 6 is determined using the marginal log-likelihood lower bound:

$$\begin{aligned} \log p(\mathbf{y}; q) &= \frac{1}{2}p(\nu + p - 1) \log(2\nu) - \frac{1}{2} \sum_{i=1}^m n_i \log(2\pi) - 4 \log(\pi) - \frac{p}{2} \log(\sigma_\beta^2) \\ &\quad - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\beta})) \right\} + \frac{1}{2}(p + L_{\text{gbl}} + m(p + L_{\text{grp}})) \\ &\quad - \frac{1}{2} \sum_{i=1}^m \log \left| \boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)}(\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} + \text{blockdiag}(M_{q(\boldsymbol{\Sigma}^{-1})}, \boldsymbol{\mu}_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}}) \right| \\ &\quad - \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}})| - \log(C_{p, \nu+p-1}) + \log(C_{p, \nu+m+p-1}) \\ &\quad + \log \Gamma \left\{ \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \right\} - \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \log(B_{q(\sigma_\varepsilon^2)}) + \log \Gamma \left(\frac{1}{2} (L_{\text{gbl}} + 1) \right) \\ &\quad - \frac{1}{2} (L_{\text{gbl}} + 1) \log(B_{q(\sigma_{\text{gbl}}^2)}) + \log \Gamma \left\{ \frac{1}{2} (m \times L_{\text{grp}} + 1) \right\} + p \log \Gamma \left(\frac{1}{2} (\nu + p) \right) \\ &\quad - \frac{1}{2} (mL_{\text{grp}} + 1) \log(B_{q(\sigma_{\text{grp}}^2)}) - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/\sigma_\varepsilon^2)} \mu_{q(1/a_\varepsilon)} \\ &\quad - \log(A_{\text{gbl}}) - \log(B_{q(a_{\text{gbl}})}) + \mu_{q(1/\sigma_{\text{gbl}}^2)} \mu_{q(1/a_{\text{gbl}})} - \log(A_{\text{grp}}) - \log(B_{q(a_{\text{grp}})}) \\ &\quad + \mu_{q(1/\sigma_{\text{grp}}^2)} \mu_{q(1/\sigma_{\text{grp}}^2)} - \sum_{j=1}^p \log(A_{\Sigma, j}) - \frac{1}{2} (\nu + p - 1) \sum_{j=1}^p \log(B_{q(a_{\Sigma, j})}) \\ &\quad + \sum_{j=1}^p \nu \left(M_{q(\boldsymbol{\Sigma}^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma, j})} - \frac{1}{2} (\nu + p + m - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma})}|. \end{aligned}$$

Set up initial values:

$\mu_{q(1/\sigma_\varepsilon^2)} > 0, \mu_{q(1/a_\varepsilon)} > 0, \mu_{q(1/\sigma_{\text{gbl}}^2)} > 0, \mu_{q(1/\sigma_{\text{grp}}^2)} > 0, \mu_{q(1/a_{\Sigma, j})} > 0, 1 \leq j \leq 2, \mathbf{M}_{q(\Sigma^{-1})}$
positive definite.

Cycle through:

$\mathbf{S} \leftarrow \mathbf{0}; \quad \mathbf{s} \leftarrow \mathbf{0}$

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

$$\mathbf{G}_i \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}_i^{\text{gbl}})^\top \mathbf{Z}_i^{\text{R}}$$

$$\mathbf{H}_i \leftarrow \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} + \text{blockdiag} \left(\mathbf{M}_{q(\Sigma^{-1})}, \mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}} \right) \right\}^{-1}$$

$$\mathbf{S} \leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top; \quad \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{Z}_i^{\text{R}})^\top \mathbf{y}_i$$

$$\Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})} \leftarrow \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{C}^{\text{gbl}})^\top \mathbf{C}^{\text{gbl}} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_{\text{gbl}}^2)} \mathbf{I}_{L_{\text{gbl}}} \end{bmatrix} - \mathbf{S} \right\}^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{gbl}})} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})} \{ (\mathbf{C}^{\text{gbl}})^\top \mathbf{y} - \mathbf{s} \}$$

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

$$\Sigma_{q(\mathbf{u}_i^{\text{R}})} \leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})} \mathbf{G}_i \mathbf{H}_i$$

$$\boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})} \leftarrow \mathbf{H}_i \left\{ \mu_{q(1/\sigma_\varepsilon^2)} (\mathbf{Z}_i^{\text{R}})^\top \mathbf{y}_i - \mathbf{G}_i^\top \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{gbl}})} \right\}$$

$$B_{q(\sigma_\varepsilon^2)} \leftarrow \mu_{q(1/a_\varepsilon)} + \frac{1}{2} \left\| \left\| \mathbf{y} - \mathbf{C}^{\text{gbl}} \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{gbl}})} - \begin{bmatrix} \mathbf{Z}_1^{\text{R}} \boldsymbol{\mu}_{q(\mathbf{u}_1^{\text{R}})} \\ \vdots \\ \mathbf{Z}_m^{\text{R}} \boldsymbol{\mu}_{q(\mathbf{u}_m^{\text{R}})} \end{bmatrix} \right\| \right\|^2$$

$$+ \text{tr} \{ (\mathbf{C}^{\text{gbl}})^\top \mathbf{C}^{\text{gbl}} \Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})} \} + \sum_{i=1}^m \text{tr} \{ (\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} \Sigma_{q(\mathbf{u}_i^{\text{R}})} \}$$

$$- 2 \mu_{q(1/\sigma_\varepsilon^2)}^{-1} \sum_{i=1}^m \text{tr} \left(\mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})} \right) \Big]$$

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

$$B_{q(\sigma_{\text{gbl}}^2)} \leftarrow \frac{1}{2} \left\{ \left\| \boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})} \right\|^2 + \text{tr} \left(\Sigma_{q(\mathbf{u}^{\text{gbl}})} \right) \right\} + \mu_{q(1/a_{\text{gbl}})}$$

$$\mu_{q(1/\sigma_{\text{gbl}}^2)} \leftarrow \frac{1}{2} (L_{\text{gbl}} + 1) / B_{q(\sigma_{\text{gbl}}^2)}, \quad \mu_{q(1/a_{\text{gbl}})} \leftarrow 1 / \left\{ \mu_{q(1/\sigma_{\text{gbl}}^2)} + A_{\text{gbl}}^{-2} \right\}$$

$$B_{q(\sigma_{\text{grp}}^2)} \leftarrow \frac{1}{2} \left\{ \left\| \boldsymbol{\mu}_{q(\mathbf{u}^{\text{grp}})} \right\|^2 + \text{tr} \left(\Sigma_{q(\mathbf{u}^{\text{grp}})} \right) \right\} + \mu_{q(1/a_{\text{grp}})}$$

$$\mu_{q(1/\sigma_{\text{grp}}^2)} \leftarrow \frac{1}{2} (m L_{\text{grp}} + 1) / B_{q(\sigma_{\text{grp}}^2)}; \quad \mu_{q(1/a_{\text{grp}})} \leftarrow 1 / \left\{ \mu_{q(1/\sigma_{\text{grp}}^2)} + A_{\text{grp}}^{-2} \right\},$$

5.2. TWO-LEVEL GAUSSIAN RESPONSE MODEL

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

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

For $j = 1, 2$:

$$B_{q(a_{\Sigma,j})} \leftarrow \nu \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right)_{jj} + 1/A_{\Sigma,j}^2, \quad \mu_{q(1/a_{\Sigma,j})} \leftarrow \left(\frac{\nu}{2} + 1 \right) / B_{q(a_{\Sigma,j})}$$

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

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

$$\boldsymbol{\Lambda}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}, \mathbf{u}_i^{\text{R}}) \equiv E_q \left[\left\{ \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^{\text{gbl}} \end{bmatrix} - \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}})} \right\} \left\{ \mathbf{u}_i^{\text{R}} - \boldsymbol{\mu}_q(\mathbf{u}_i^{\text{R}}) \right\}^\top \right] \leftarrow -\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}})} \mathbf{G}_i \mathbf{H}_i$$

Algorithm 6: *Streamlined MFVB algorithm for the two-level Gaussian response model given in (5.8).*

The difference between this lower bound expression and the one given in (5.2.3) is replacement of the term

$$\log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|$$

with $-\sum_{i=1}^m \log \left| \mu_{q(1/\sigma_i^2)} (\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} + \text{blockdiag} \left(M_{q(\boldsymbol{\Sigma}^{-1})}, \mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}} \right) \right| - \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}})}|$,

which is used to streamline the calculation of $\log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|$. Details on this calculation are given in the appendix of Lee & Wand (2015).

In order to produce variability estimates when plotting subject-specific mean estimates, Lee & Wand (2015) propose the following partition of the q -density covariance matrix of the vector of coefficients for the i th subject:

$$\text{COV}_q \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u}^{\text{gbl}} \\ \mathbf{u}_i^{\text{R}} \end{bmatrix} \right) = \begin{bmatrix} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}})} & \boldsymbol{\Lambda}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}, \mathbf{u}_i^{\text{R}}) \\ \boldsymbol{\Lambda}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}, \mathbf{u}_i^{\text{R}})^\top & \boldsymbol{\Sigma}_{q(\mathbf{u}_i^{\text{R}})} \end{bmatrix}.$$

It is important to keep in mind that the time savings when using a streamlined MFVB approach for the Bayesian hierarchical two-level Gaussian response model given in (5.19) will not be as impressive as the ones achieved by using the two-level model in Lee & Wand (2015). The reason for this lies in the structure of the \mathbf{H}_i matrix. For example, in Lee & Wand (2015) there exists no zero's in \mathbf{H}_i , however in (5.19), since we incorporate the variance of the spline components for the group specific curve of each individual there exists

blocks of zeros concerning the $\mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}}$ term in Algorithm 6. As a result, this should not be considered a fully streamlined version of Algorithm 5 and further investigation into a more robust version of Algorithm 6 is certainly warranted.

5.2.5 Simulation study

We carried out an extensive simulation study in order to assess the speed of Algorithm 6 against that of Algorithm 5. We generated 25 data-sets according to (5.1) where the $x_{ij} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1)$, $u_{\ell}^{\text{gbl}} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\text{gbl}}^2)$, $\varepsilon_{ij} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\varepsilon}^2)$, $\boldsymbol{\delta}_i \stackrel{\text{ind.}}{\sim} \text{N}(\mathbf{0}, \boldsymbol{\Sigma})$ and the $u_{i\ell}^{\text{grp}} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\text{grp}}^2)$. The true parameter values were specified as

$$\beta_0 = 0.2, \quad \beta_1 = 1.8, \quad \sigma_{\text{gbl}}^2 = 1.5, \quad \sigma_{\varepsilon}^2 = 0.05, \quad \boldsymbol{\Sigma} = \begin{bmatrix} 0.4 & 0.1 \\ 0.1 & 0.4 \end{bmatrix} \quad \text{and} \quad \sigma_{\text{grp}}^2 = 0.29.$$

The number of subjects between simulation studies was varied, where $m \in \{25, 75, 125, 175\}$ and the within subject sample size remained constant at $n_i = 10$, $1 \leq i \leq m$.

Algorithm 5 and 6 were implemented in the R programming language. These computations were performed on a laptop computer (Mac OS X; 2.8 GHz processor, 16 GBytes of random access memory). Table 5.1 summarises the computation times for each approach. As the number of subjects m increase moderately, the average computing time for the naïve approach increases very quickly. The average time increase for the naïve approach from $m = 25$ to $m = 175$ is approximately 3 hours, compared to the streamlined approach which was just under 1 minute. This is represented in the *Ratio* column of Table 5.1 where we can see exactly how much faster the streamlined approach is in comparison to the naïve approach. Keeping in mind that the largest value of m is 175, the possibilities of time saving given larger values of m seem endless.

m	Naïve	Streamlined	Ratio
25	37.16 (17.49)	9.52 (3.75)	3.86
75	761.16 (248.39)	20.44 (6.13)	37.24
125	3178.27 (562.68)	31.45 (5.22)	101.06
175	11181.61 (4587.89)	56.49 (22.36)	197.94

Table 5.1: Average (standard deviation) run time in seconds for naïve and streamlined MFVB fitting of the two-level Gaussian response model in (5.1).

5.2.6 Application to data from a growth study

In this section we provide an example of streamlined MFVB fitting for growth data on children in Indianapolis (source: Pratt *et al.*, 1989). These data possess a two-level structure where the repeated measurements on height and age are at the first level, and the persons being measured are at the second level. In this data-set a person is identified as being either *black* or *white* by the use of the indicator variable `black`. We only use the information from this data-set relating to male individuals. This gives $m = 116$, and each n_i takes on a value between 10 and 26 repeated measurements each. The variables and their description are given in Table 5.2.

Variable	Description
<code>id</code>	person identifier
<code>height</code>	person's height in centimeters
<code>age</code>	person's age
<code>black</code>	indicator of person being black

Table 5.2: *Description of the Indiana growth data for male subjects.*

We are interested in what the effect of being either *black* or *white* would have on a person's height. Thus an interaction term would be necessary in our model to account for this. The Indiana growth data is divided into two groups:

$$A \equiv \text{white males} \quad B \equiv \text{black males},$$

and to quantify the difference between the two groups we extend (5.2) to:

$$\begin{aligned}
 f_A(\mathbf{age}) &= \beta_0^A + \beta_1^A \mathbf{age} + \sum_{\ell=1}^{L^{\text{gbl}}} u_{A,\ell}^{\text{gbl}} z_{\ell}^{\text{gbl}}(\mathbf{age}), \\
 f_B(\mathbf{age}) &= \beta_0^A + \beta_0^{\text{BvsA}} + (\beta_1^A + \beta_1^{\text{BvsA}}) \mathbf{age} + \sum_{\ell=1}^{L^{\text{gbl}}} u_{B,\ell}^{\text{gbl}} z_{\ell}^{\text{gbl}}(\mathbf{age}), \\
 \text{and } g_i(\mathbf{age}) &= \delta_{0i} + \delta_{1i} \mathbf{age} + \sum_{\ell=1}^{L^{\text{grp}}} u_{i\ell}^{\text{grp}} z_{\ell}^{\text{grp}}(\mathbf{age}).
 \end{aligned} \tag{5.14}$$

This allows one to estimate the contrast function given in (5.16). We define the parameter vectors to be

$$\boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0^A \\ \beta_1^A \\ \beta_0^{\text{BvsA}} \\ \beta_1^{\text{BvsA}} \end{bmatrix}, \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}_A^{\text{gbl}} \\ \mathbf{u}_B^{\text{gbl}} \\ \mathbf{u}^{\text{R}} \end{bmatrix}$$

where

$$\mathbf{u}_A^{\text{gbl}} \equiv \begin{bmatrix} u_{A,1}^{\text{gbl}} \\ \vdots \\ u_{A,L^{\text{gbl}}}^{\text{gbl}} \end{bmatrix}, \quad \mathbf{u}_B^{\text{gbl}} \equiv \begin{bmatrix} u_{B,1}^{\text{gbl}} \\ \vdots \\ u_{B,L^{\text{gbl}}}^{\text{gbl}} \end{bmatrix}.$$

In addition, the design matrices \mathbf{X} and \mathbf{Z} are

$$\mathbf{X} \equiv \begin{bmatrix} \mathbf{1} & \text{age}_1 & \mathbf{I}_1^A \odot \mathbf{1} & \mathbf{I}_1^A \odot \text{age}_1 \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{1} & \text{age}_m & \mathbf{I}_m^A \odot \mathbf{1} & \mathbf{I}_m^A \odot \text{age}_m \end{bmatrix}, \quad \mathbf{Z} \equiv [\mathbf{Z}_A^{\text{gbl}} \mid \mathbf{Z}_B^{\text{gbl}} \mid \mathbf{Z}^R],$$

where

$$\mathbf{Z}_A^{\text{gbl}} \equiv \begin{bmatrix} \mathbf{I}_1^A \odot z_1^{\text{gbl}}(\text{age}_1) & \dots & \mathbf{I}_1^A \odot z_{L^{\text{gbl}}}^{\text{gbl}}(\text{age}_1) \\ \vdots & \ddots & \vdots \\ \mathbf{I}_m^A \odot z_1^{\text{gbl}}(\text{age}_m) & \dots & \mathbf{I}_m^A \odot z_{L^{\text{gbl}}}^{\text{gbl}}(\text{age}_m) \end{bmatrix}$$

and

$$I_{ij}^A \equiv \begin{cases} 1 & \text{if } (\text{age}_{ij}, \text{height}_{ij}) \text{ is of type A,} \\ 0 & \text{if } (\text{age}_{ij}, \text{height}_{ij}) \text{ is of type B.} \end{cases}$$

The $n_i \times 1$ vector \mathbf{I}_i^A consists of the I_{ij}^A and the design matrix $\mathbf{Z}_B^{\text{gbl}}$ is defined in a similar way to $\mathbf{Z}_A^{\text{gbl}}$ but with \mathbf{I}_i^A replaced by $\mathbf{1} - \mathbf{I}_i^A$. The full Bayesian model is then

$$\begin{aligned} \text{height} \mid \beta, \mathbf{u}, \sigma_\varepsilon^2 &\sim \text{N}(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), & \beta &\sim \text{N}(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_2), \\ \mathbf{u} \mid (\sigma_{\text{gbl}}^A)^2, (\sigma_{\text{gbl}}^B)^2, \Sigma, \sigma_{\text{grp}}^2 &\sim \text{N}\left(\mathbf{0}, \begin{bmatrix} (\sigma_{\text{gbl}}^A)^2 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\sigma_{\text{gbl}}^B)^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag}(\Sigma, \sigma_{\text{grp}}^2 \mathbf{I}_{L_{\text{grp}}}) \end{bmatrix}\right), \\ \sigma_\varepsilon &\sim \text{Half-Cauchy}(A_\varepsilon), & \sigma_{\text{gbl}}^A &\sim \text{Half-Cauchy}(A_{\text{gbl}}^A), \\ \sigma_{\text{gbl}}^B &\sim \text{Half-Cauchy}(A_{\text{gbl}}^B), & \sigma_{\text{grp}} &\sim \text{Half-Cauchy}(A_{\text{grp}}). \\ \Sigma \mid a_{\Sigma,1}, a_{\Sigma,2} &\sim \text{Inverse-Wishart}(\nu + 2 - 1, 2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})), \\ a_{\Sigma,j} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\Sigma,j}^2\right), & 1 \leq j &\leq 2. \end{aligned} \tag{5.15}$$

We use an equivalent set-up and partition of the mixed model framework as shown in Sections 5.2.1 and 5.2.4. Model (5.15) was fit using the streamlined MFVB algorithm as shown in Algorithm 6, where the variables `height` and `age` were standardised to have zero mean and unit variance and the values of the hyperparameters were all set to 10^5 . The `rstan` package which uses `Stan` code, was used for performing MCMC with a burn-in of 5000, number of further iterations of 5000 and a thinning factor of 5. The MFVB iterations were terminated when the increase in the corresponding $\log p(\text{height}; q)$ fell below 10^{-8} .

As mentioned in Chapter 2, making time comparisons fair between MFVB and MCMC is quite difficult, due to convergence for each approximation method being vastly different.

5.2. TWO-LEVEL GAUSSIAN RESPONSE MODEL

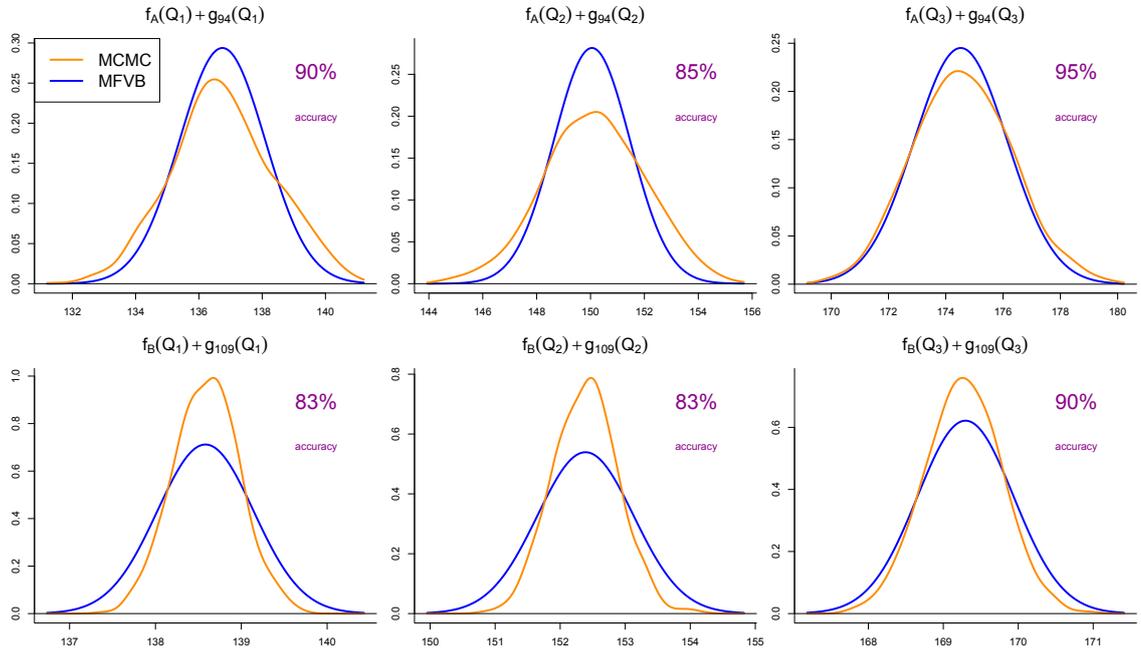


Figure 5.3: *Approximate posterior density functions obtained via MCMC and streamlined MFVB for the growth Indiana data set. MFVB accuracy scores are displayed.*

Here we make use of **Stan** as opposed to **BUGS**, which we have been using until now. **Stan** is known to generally converge faster than **BUGS** as it is based on Hamiltonian Monte Carlo sampling. **Stan** is also preferred for certain problems with a complex structure, such as multilevel models.

Figure 5.3 illustrates the accuracy of the streamlined MFVB approach against the MCMC benchmark. The parameters being monitored are the quartiles of the curve estimates of the 94th and 109th individual. These are represented as $f_A(Q_i) + g_{94}(Q_i)$ and $f_B(Q_i) + g_{109}(Q_i)$ for $i = 1, 2, 3$, where the 94th individual is *white* and the 109th individual is *black*. The accuracy of the curve estimates for these parameters are reasonably high considering the complexity of the model.

Figure 5.4 provides illustration of the approximate inference achieved by both streamlined MFVB and MCMC for the two-level Gaussian response model. It appears that the subject-specific curve estimates and pointwise 95% credible sets are almost indistinguishable for both approaches, implying that the streamlined MFVB approach incurs little loss of accuracy. It is immediately obvious that the streamlined MFVB fits are agreeing very well with the longitudinal data for each individual. The MCMC fit took almost 3 hours, whilst the streamlined MFVB fit took only 4 seconds.

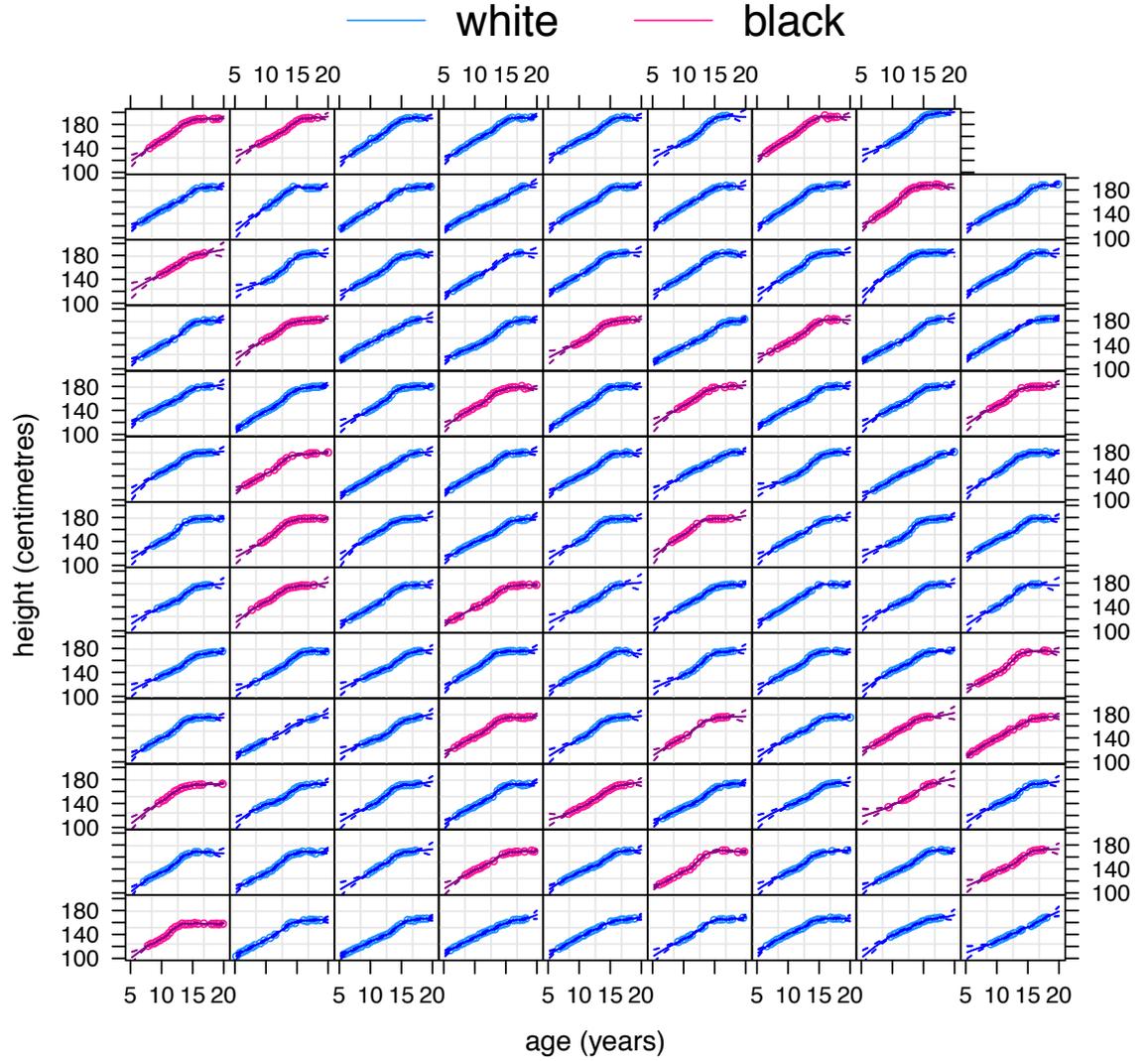


Figure 5.4: Group-specific estimates and pointwise 95% credible sets obtained via streamlined MFVB and MCMC for the Indiana growth data-set for male subjects.

Figure 5.5 illustrates the effect of the variable `black` on the height of male individuals. This is given by the *contrast curve*

$$\begin{aligned}
 c(\text{age}) &\equiv f_B(\text{age}) - f_A(\text{age}) \\
 &= \beta_0^{\text{BvsA}} + \beta_1^{\text{BvsA}} \text{age} + \sum_{\ell=1}^{L_{\text{gbl}}} (u_{\text{B},\ell}^{\text{gbl}} - u_{\ell,\text{W}}^{\text{gbl}}) z_{\ell}^{\text{gbl}}(\text{age}).
 \end{aligned} \tag{5.16}$$

A contrast curve was generated for both MCMC and streamlined MFVB. It is evident that at age 5, on average, *black* individuals are approximately 3.5cm taller than *white* individuals. However, as they reach adolescence, *black* individuals are, on average, taller by approximately 6cm. Finally, as males reach their early 20's, the height differences decrease to zero. In addition, the agreement of the curves produced by MCMC and

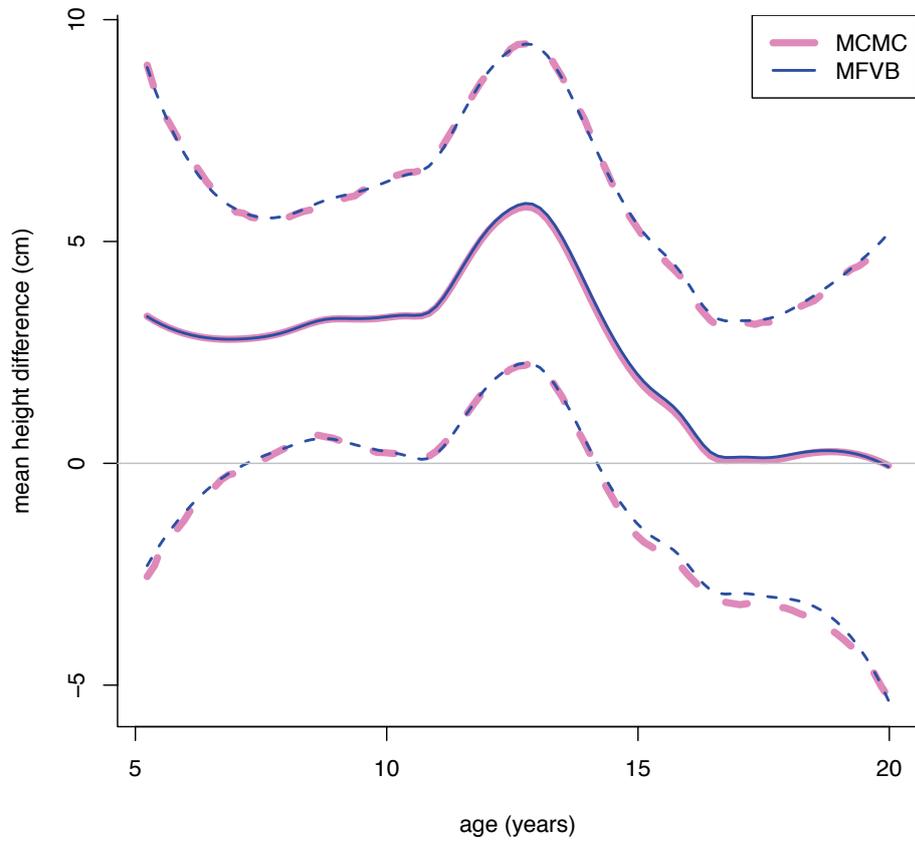


Figure 5.5: *Approximate contrast curves obtained via MCMC and streamlined MFVB for the mean height difference between black and white males.*

streamlined MFVB are quite pleasing, indicating little loss of accuracy.

5.3 Three-level Gaussian response model

In this section, we extend the methodology shown in the previous section to cater to the three-level Gaussian response model for both the naïve and streamlined MFVB situations. We apply this methodology to a data set which exhibits a three-level longitudinal structure. These data represent frequency dependent backscatter coefficients for induced tumors in rodents (source: Simpson, 2013). There are several *slices* or probe locations in which each tumor is scanned. Five different transducers were used on various slices to perform scanning. The variables and their descriptions are given in Table 5.3. Figure 5.6

Variable	Description
<code>freq</code>	frequency measurement
<code>bsc</code>	backscatter coefficient (converted to decibel scale)
<code>L9h4</code>	indicator of transducer L9-4 being used
<code>L40</code>	indicator of transducer L40 being used
<code>MS200</code>	indicator of transducer MS200 being used
<code>MS400</code>	indicator of transducer MS400 being used

Table 5.3: *Description of the rodent tumor dataset.*

illustrates the information from 10 rodents given in this dataset. Each panel represents a different tumor corresponding to a rodent. The x -axis represents the frequency and the y -axis represents the backscatter coefficients for each tumor. Each curve in each panel corresponds to a different slice. Furthermore, within each panel, the colours correspond to different transducers being used according to:

green L14-5
red L9-4
blue L40
orange MS200
purple MS400.

We use similar notation to that used in Section 5.2, where for example

$\text{bsc}_{ijk} \equiv k$ th backscatter coefficient for the j th slice within the i th tumor.

The generic three-level Gaussian response model extension, specific to the data presented in Figure 5.6 is

$$\text{bsc}_{ijk} \sim N\left(f(\text{freq}_{ijk}) + g_i(\text{freq}_{ijk}) + h_{ij}(\text{freq}_{ijk}), \sigma_\varepsilon^2\right), \quad (5.17)$$

$$1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij},$$

5.3. THREE-LEVEL GAUSSIAN RESPONSE MODEL

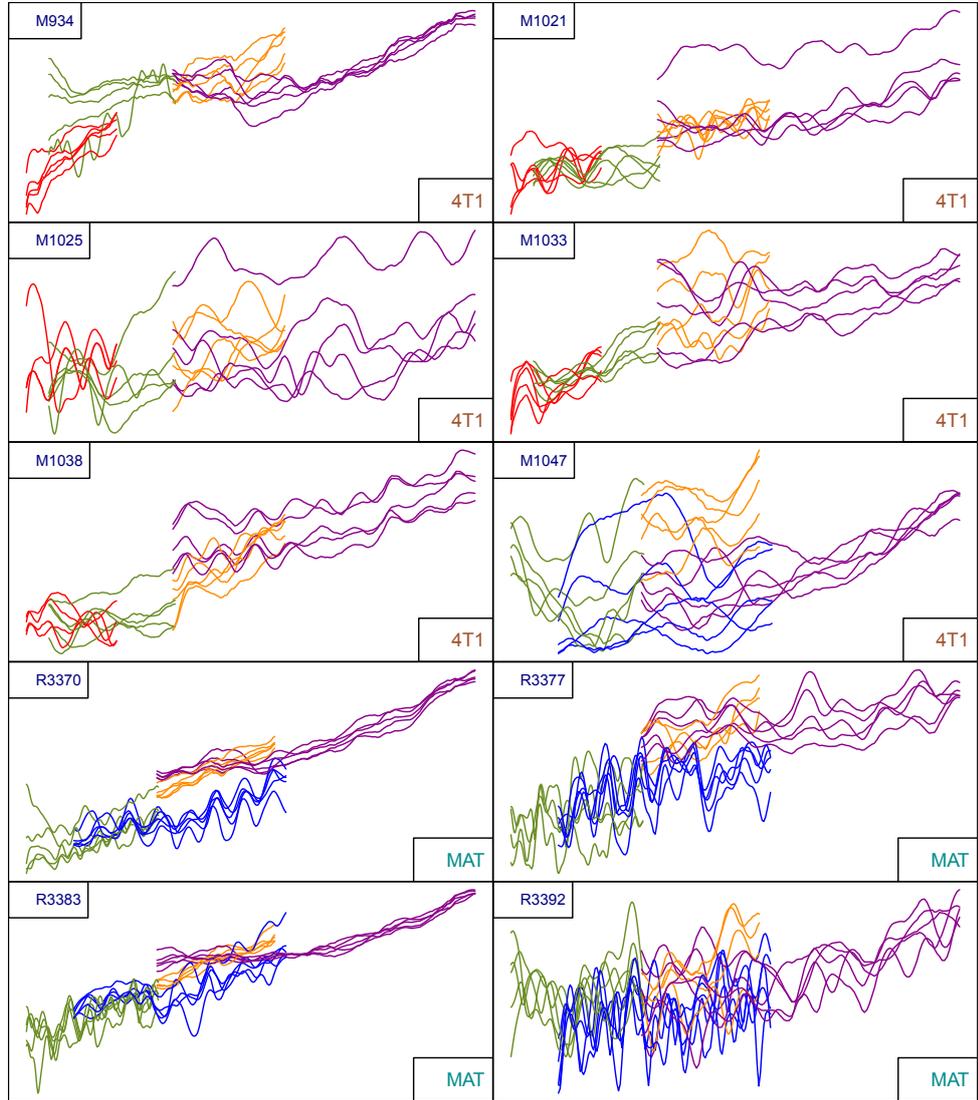


Figure 5.6: *Rodent tumor data set showing frequency dependent backscatter coefficients for two different types of induced tumors in rodents. Each panel is for a different tumor (name shown in top left corner). The type of tumor is shown in the bottom right corner of each panel.*

where we have $m = 10$ tumors, the n_i range from being either 18, 19 or 20 slices (after removing outlier slices) and the $o_{ij} = 128$ is constant throughout and we use o from now on to represent the measurements in each slice. The global mean function has the form

$$f(x) = \beta_0 + \beta_1 x + \beta_2 \text{L9h4} + \beta_3 \text{L40} + \beta_4 \text{MS200} + \beta_5 \text{MS400} + \sum_{\ell=1}^{L_{\text{gbl}}} u_{\ell}^{\text{gbl}} z_{\ell}^{\text{gbl}}(x),$$

5.3. THREE-LEVEL GAUSSIAN RESPONSE MODEL

where $u_\ell^{\text{gbl}} | \sigma_{\text{gbl}}^2 \stackrel{\text{ind.}}{\sim} N\{0, \sigma_{\text{gbl}}^2\}$. In order to model the deviation from f for the i th tumor, we further extend our model to include the function

$$g_i(x) = \delta_{0i} + \delta_{1i}x + \sum_{\ell=1}^{L_{\text{grpout}}} u_\ell^{\text{grpout}} z_\ell^{\text{grpout}}(x), \quad \begin{bmatrix} \delta_{0i} \\ \delta_{1i} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}_{\text{grpout}}),$$

$$u_\ell^{\text{grpout}} | \sigma_{\text{grpout}}^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_{\text{grpout}}^2),$$

where ‘grpout’ refers to each rodent’s information. In addition, to model the deviation from g_i for the ij th slice, we also make use of the function

$$h_{ij}(x) = \gamma_{0ij} + \gamma_{1ij}x + \sum_{\ell=1}^{L_{\text{grpinn}}} u_\ell^{\text{grpinn}} z_\ell^{\text{grpinn}}(x), \quad \begin{bmatrix} \gamma_{0ij} \\ \gamma_{1ij} \end{bmatrix} \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma}_{\text{grpinn}}),$$

$$u_\ell^{\text{grpinn}} | \sigma_{\text{grpinn}}^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_{\text{grpinn}}^2),$$

where ‘grpinn’ refers to the information involving the measurements for each slice. In order to express our model in a mixed model framework, we first define the following vectors and matrices:

$$\boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{bmatrix}, \quad \mathbf{u}^{\text{gbl}} \equiv \begin{bmatrix} u_1^{\text{gbl}} \\ \vdots \\ u_{L_{\text{gbl}}}^{\text{gbl}} \end{bmatrix}, \quad \mathbf{u}_{\text{full}}^{\text{grpout}} \equiv \begin{bmatrix} \boldsymbol{\delta}_1 \\ \mathbf{u}_1^{\text{grpout}} \\ \vdots \\ \boldsymbol{\delta}_m \\ \mathbf{u}_m^{\text{grpout}} \end{bmatrix}, \quad \mathbf{u}_{\text{full}}^{\text{grpinn}} \equiv \begin{bmatrix} \gamma_{11} \\ \mathbf{u}_{11}^{\text{grpinn}} \\ \vdots \\ \gamma_{1n_1} \\ \mathbf{u}_{1n_1}^{\text{grpinn}} \\ \vdots \\ \gamma_{m1} \\ \mathbf{u}_{m1}^{\text{grpinn}} \\ \vdots \\ \gamma_{mn_m} \\ \mathbf{u}_{mn_m}^{\text{grpinn}} \end{bmatrix},$$

$$\mathbf{u} \equiv \begin{bmatrix} \mathbf{u}^{\text{gbl}} \\ \mathbf{u}_{\text{full}}^{\text{grpout}} \\ \mathbf{u}_{\text{full}}^{\text{grpinn}} \end{bmatrix},$$

where the $\boldsymbol{\delta}_i$ and $\mathbf{u}_i^{\text{grpout}}$ are defined as in section 5.2.1. We define $\boldsymbol{\gamma}_{ij} \equiv [\gamma_{0ij} \ \gamma_{1ij}]^\top$ and $\mathbf{u}_{ij}^{\text{grpinn}} \equiv [u_1^{\text{grpinn}}, \dots, u_{L_{\text{grpinn}}}^{\text{grpinn}}]^\top$. The design matrices \mathbf{X} and \mathbf{Z} are

$$\mathbf{X} \equiv \begin{bmatrix} \mathbf{1} & \text{freq}_1 & \text{L9h4}_1 & \text{L40}_1 & \text{MS200}_1 & \text{MS400}_1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{1} & \text{freq}_m & \text{L9h4}_m & \text{L40}_m & \text{MS200}_m & \text{MS400}_m \end{bmatrix},$$

$$\mathbf{Z} \equiv \left[\mathbf{Z}^{\text{gbl}} \mid \text{blockdiag} \left(\mathbf{X}_i \mid \mathbf{Z}_i^{\text{grpout}} \right) \mid \text{blockdiag} \left(\text{blockdiag} \left(\mathbf{X}_{ij} \mid \mathbf{Z}_{ij}^{\text{grpinn}} \right) \right) \right],$$

where

$$\mathbf{Z}^{\text{gbl}} \equiv \begin{bmatrix} z_1^{\text{gbl}}(\text{freq}_1) & \dots & z_{L_{\text{gbl}}}^{\text{gbl}}(\text{freq}_1) \\ \vdots & \ddots & \vdots \\ z_1^{\text{gbl}}(\text{freq}_m) & \dots & z_{L_{\text{gbl}}}^{\text{gbl}}(\text{freq}_m) \end{bmatrix},$$

$$\mathbf{Z}_i^{\text{grpout}} \equiv \begin{bmatrix} z_1^{\text{grpout}}(\text{freq}_{i1}) & \dots & z_{L_{\text{grpout}}}^{\text{grpout}}(\text{freq}_{i1}) \\ \vdots & \ddots & \vdots \\ z_1^{\text{grpout}}(\text{freq}_{in_i}) & \dots & z_{L_{\text{grpout}}}^{\text{grpout}}(\text{freq}_{in_i}) \end{bmatrix},$$

$$\mathbf{Z}_{ij}^{\text{grpinn}} \equiv \begin{bmatrix} z_1^{\text{grpinn}}(\text{freq}_{ij1}) & \dots & z_{L_{\text{grpinn}}}^{\text{grpinn}}(\text{freq}_{ij1}) \\ \vdots & \ddots & \vdots \\ z_1^{\text{grpinn}}(\text{freq}_{ijo}) & \dots & z_{L_{\text{grpinn}}}^{\text{grpinn}}(\text{freq}_{ijo}) \end{bmatrix},$$

with freq_i being the $n_i \times 1$ vector containing the freq_{ijk} and \mathbf{X}_i is the subset of \mathbf{X} corresponding to the i th tumor. We can now represent our model using the mixed model framework:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad (5.18)$$

$$\mathbf{G} \equiv \text{Cov}(\mathbf{u}) =$$

$$\begin{bmatrix} \sigma_{\text{gbl}}^2 \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left(\boldsymbol{\Sigma}_{\text{grpout}}, \sigma_{\text{grpout}}^2 \mathbf{I}_{L_{\text{grpout}}} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag} \left(\boldsymbol{\Sigma}_{\text{grpinn}}, \sigma_{\text{grpinn}}^2 \mathbf{I}_{L_{\text{grpinn}}} \right) \end{bmatrix}$$

5.3.1 Bayesian three-level Gaussian response model

Much like in the two-level case, fitting (5.18) using standard mixed model software such as `lme()` in R is possible, but extremely time intensive. We opt to work with a Bayesian version of (5.18) which allows direct implementation in standard Bayesian software. We enforce $\sigma_{\boldsymbol{\varepsilon}}$, σ_{gbl} , $\boldsymbol{\Sigma}_{\text{grpout}}$, σ_{grpout} , $\boldsymbol{\Sigma}_{\text{grpinn}}$ and σ_{grpinn} to possess half-t prior distributions on $(-1, 1)$, as suggested by Huang *et al.* (2013). The standard deviation parameter priors are achieved through auxiliary variable constructions as exhibited in Result 1.4.6. The full

Bayesian three-level Gaussian response model for the rodent tumor data 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(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_6), & \mathbf{u} | \mathbf{G} &\sim N(\mathbf{0}, \mathbf{G}), \\
 \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), \\
 \sigma_{\text{gbl}}^2 | a_{\text{gbl}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{gbl}}\right), & a_{\text{gbl}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, A_{\text{gbl}}^{-2}\right), \\
 \boldsymbol{\Sigma}_{\text{grpout}} | a_{\Sigma_{\text{grpout},1}}, a_{\Sigma_{\text{grpout},2}} &\sim \text{Inverse-Wishart}\left(\nu_{\text{grpout}} + 2 - 1, \right. \\
 & \left. 2\nu_{\text{grp}} \text{diag}\left(1/a_{\Sigma_{\text{grpout},1}}, 1/a_{\Sigma_{\text{grpout},2}}\right)\right), \\
 a_{\Sigma_{\text{grpout},j}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\Sigma_{\text{grpout},j}}^2\right), & 1 \leq j \leq 2, \\
 \sigma_{\text{grpout}}^2 | a_{\text{grpout}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{grpout}}\right), \\
 a_{\text{grpout}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{grpout}}^2\right), \\
 \boldsymbol{\Sigma}_{\text{grpinn}} | a_{\Sigma_{\text{grpinn},1}}, a_{\Sigma_{\text{grpinn},2}} &\sim \text{Inverse-Wishart}\left(\nu_{\text{grpinn}} + 2 - 1, \right. \\
 & \left. 2\nu_{\text{grpinn}} \text{diag}\left(1/a_{\Sigma_{\text{grpinn},1}}, 1/a_{\Sigma_{\text{grpinn},2}}\right)\right), \\
 a_{\Sigma_{\text{grpinn},j}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\Sigma_{\text{grpinn},j}}^2\right), & 1 \leq j \leq 2, \\
 \sigma_{\text{grpinn}}^2 | a_{\text{grpinn}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_{\text{grpinn}}\right), \\
 a_{\text{grpinn}} &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_{\text{grpinn}}^2\right),
 \end{aligned} \tag{5.19}$$

where a_ε , a_{gbl} , $\mathbf{a}_{\Sigma_{\text{grpout}}}$, a_{grpout} , $\mathbf{a}_{\Sigma_{\text{grpinn}}}$ and a_{grpinn} are the auxiliary variables used for the scale parameters in the model. Fitting (5.19) involves the setting of hyperparameters σ_β , A_ε , A_{gbl} , ν_{grpout} , $A_{\Sigma_{\text{grpout},1}}$, $A_{\Sigma_{\text{grpout},2}}$, A_{grpout} , ν_{grpinn} , $A_{\Sigma_{\text{grpinn},1}}$, $A_{\Sigma_{\text{grpinn},2}}$ and A_{grpinn} . To ensure that Property 4 of Huang *et al.* (2013) holds, we set $\nu_{\text{grpout}} = \nu_{\text{grpinn}} = 2$. For the remaining hyperparameters, non-informativity is achieved when their values are large, e.g., 10^5 .

5.3.2 Mean field variational Bayes methodology

The essence of the MFVB approach for (5.19) is to approximate the joint posterior density function with an approximate product density form

$$\begin{aligned}
 &q(\boldsymbol{\beta}, \mathbf{u}, a_\varepsilon, a_{\text{gbl}}, \mathbf{a}_{\Sigma_{\text{grpout}}}, a_{\text{grpout}}, \mathbf{a}_{\Sigma_{\text{grpinn}}}, a_{\text{grpinn}}) \\
 &\times q\left\{\sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}_{\text{grpout}}, \sigma_{\text{grpout}}^2, \boldsymbol{\Sigma}_{\text{grpinn}}, \sigma_{\text{grpinn}}^2\right\}.
 \end{aligned}$$

Using induced factorisation theory, it is shown that a further factorisation is of the form:

$$\begin{aligned}
 &q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma_\varepsilon^2) q(\sigma_{\text{gbl}}^2) q(\boldsymbol{\Sigma}_{\text{grpout}}) q(\sigma_{\text{grpout}}^2) q(a_\varepsilon) q(a_{\text{gbl}}) \\
 &\times \left\{ \prod_{j=1}^2 q(a_{\Sigma_{\text{grpout},j}}) \right\} q(a_{\text{grpout}}) \left\{ \prod_{j=1}^2 q(a_{\Sigma_{\text{grpinn},j}}) \right\} q(a_{\text{grpinn}}).
 \end{aligned} \tag{5.20}$$

As mentioned in section 5.2.3 the optimal q -densities can be obtained through expressions such as that given in (5.11). The full conditional distributions that we deal with all

5.3. THREE-LEVEL GAUSSIAN RESPONSE MODEL

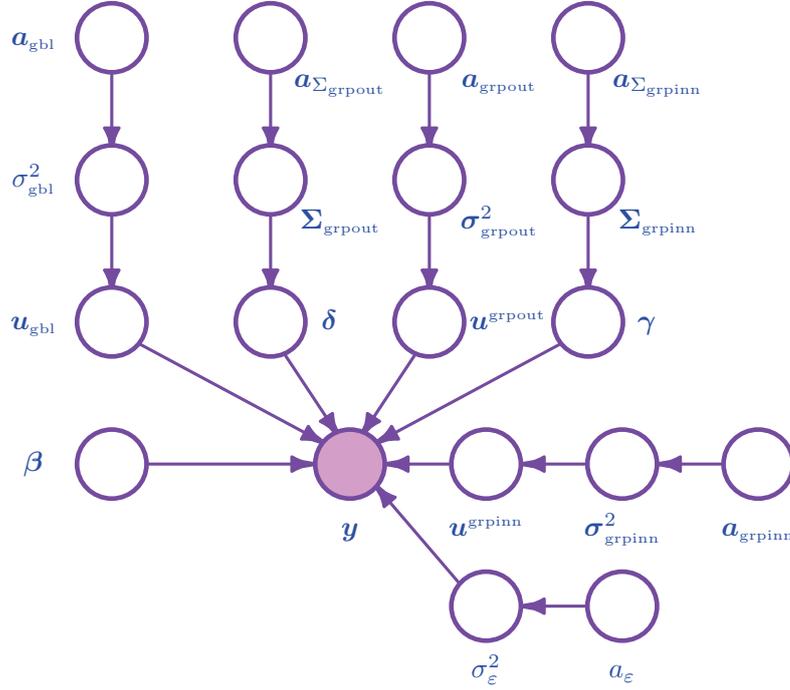


Figure 5.7: DAG corresponding to the Bayesian hierarchical three-level Gaussian response model in (5.19).

have closed form expressions, resulting in our optimal q -densities possessing closed form expressions. In light of this, the factors given in (5.20) have the following forms:

$$\begin{aligned}
 q^*(\boldsymbol{\beta}, \mathbf{u}) &\sim N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \\
 q^*(\sigma_{\epsilon}^2) &\sim \text{Inverse-Gamma}\left(\frac{1}{2}\left(\sum_{i=1}^m n_i + 1\right), B_q(\sigma_{\epsilon}^2)\right) \\
 q^*(\sigma_{\text{gbl}}^2) &\sim \text{Inverse-Gamma}\left(\frac{1}{2}(L_{\text{gbl}} + 1), B_q(\sigma_{\text{gbl}}^2)\right) \\
 q^*(\boldsymbol{\Sigma}_{\text{grpout}}) &\sim \text{Inverse-Wishart}(\nu_{\text{grpout}} + m + 1, B_q(\boldsymbol{\Sigma}_{\text{grpout}})) \\
 q^*(\sigma_{\text{grpout}}^2) &\sim \text{Inverse-Gamma}\left(\frac{1}{2}(m L_{\text{grpout}} + 1), B_q(\sigma_{\text{grpout}}^2)\right) \\
 q^*(\boldsymbol{\Sigma}_{\text{grpinn}}) &\sim \text{Inverse-Wishart}(\nu_{\text{grpinn}} + \sum_{i=1}^n n_i + 1, B_q(\boldsymbol{\Sigma}_{\text{grpinn}})), \\
 q^*(\sigma_{\text{grpinn}}^2) &\sim \text{Inverse-Gamma}\left\{\frac{1}{2}\left(L_{\text{grpinn}} \sum_{i=1}^m n_i + 1\right), B_q(\sigma_{\text{grpinn}}^2)\right\}, \\
 q^*(a_{\epsilon}) &\sim \text{Inverse-Gamma}(1, B_q(a_{\epsilon})), \quad q^*(a_{\text{gbl}}) \sim \text{Inverse-Gamma}(1, B_q(a_{\text{gbl}})) \\
 q^*(a_{\text{grpout}}) &\sim \text{Inverse-Gamma}(1, B_q(a_{\text{grpout}})), \\
 q^*(a_{\text{grpinn}}) &\sim \text{Inverse-Gamma}(1, B_q(a_{\text{grpinn}})), \\
 q^*(a_{\Sigma_{\text{grpout}}, j}) &\sim \text{Inverse-Gamma}\left(\frac{\nu_{\text{grpout}}}{2} + 1, B_q(a_{\Sigma_{\text{grpout}}, j})\right) \quad 1 \leq j \leq 2, \\
 q^*(a_{\Sigma_{\text{grpinn}}, j}) &\sim \text{Inverse-Gamma}\left(\frac{\nu_{\text{grpinn}}}{2} + 1, B_q(a_{\Sigma_{\text{grpinn}}, j})\right) \quad 1 \leq j \leq 2.
 \end{aligned} \tag{5.21}$$

5.3. THREE-LEVEL GAUSSIAN RESPONSE MODEL

The derivation of (5.21) is similar to that shown for the two-level model in Appendix 5.A and the parameters in (5.21) are defined analogously to that described in Section 5.2.3. The values of these optimal parameters are obtained through an iterative scheme presented as Algorithm 7, where

$$\mathbf{M}_q(\mathbf{G}^{-1}) = \begin{bmatrix} \mu_{q(1/\sigma_{\text{gbl}}^2)} \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq i \leq m} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpout}}^{-1})}, \mu_{q(1/\sigma_{\text{grpout}}^2)} \mathbf{I}_{L_{\text{grpout}}} \right) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \text{blockdiag}_{1 \leq j \leq \sum_{i=1}^m n_i} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpinn}}^{-1})}, \mu_{q(1/\sigma_{\text{grpinn}}^2)} \mathbf{I}_{L_{\text{grpinn}}} \right) \end{bmatrix}.$$

We define

$$\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})} \equiv \text{sub-vector of } \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \text{ corresponding to } \mathbf{u}^{\text{gbl}}$$

and

$$\boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{gbl}})} \equiv \text{sub-matrix of } \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \text{ corresponding to } \mathbf{u}^{\text{gbl}}.$$

The vectors $\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grpout}})}$, $\boldsymbol{\mu}_{q(\boldsymbol{\delta}_i)}$, $\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grpinn}})}$, $\boldsymbol{\mu}_{q(\boldsymbol{\gamma}_{ij})}$ and their corresponding variance covariance matrices are defined analogously. Algorithm 7 makes use of the variational lower bound on the marginal log-likelihood. For model (5.19) and product restriction (5.20) it

Set up initial values:

$$\begin{aligned} & \mu_{q(1/\sigma_\varepsilon^2)}, \mu_{q(1/\sigma_{\text{gbl}}^2)}, \mu_{q(1/\sigma_{\text{grpout}}^2)}, \mu_{q(1/\sigma_{\text{grpinn}}^2)}, \mu_{q(1/a_\varepsilon)}, \mu_{q(1/a_{\text{gbl}})}, \mu_{q(1/a_{\text{grpout}})}, \\ & \mu_{q(1/a_{\text{grpinn}})} > 0, \mu_{q(1/a_{\Sigma_{\text{grpout},j}})}, \mu_{q(1/a_{\Sigma_{\text{grpinn},j}})} > 0, \quad 1 \leq j \leq 2, \quad \Sigma_{q(\beta, \mathbf{u})}, \\ & \mathbf{M}_{q(\Sigma_{\text{grpout}}^{-1})}, \mathbf{M}_{q(\Sigma_{\text{grpinn}}^{-1})} \text{ positive definite.} \end{aligned}$$

Cycle through:

$$\begin{aligned} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} & \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta, \mathbf{u})} \mathbf{C}^\top \mathbf{y} \\ \Sigma_{q(\beta, \mathbf{u})} & \leftarrow \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{M}_{q(\mathbf{G}^{-1})} \end{bmatrix} \right)^{-1} \\ B_{q(\sigma_\varepsilon^2)} & \leftarrow \frac{1}{2} \{ \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C} \Sigma_{q(\beta, \mathbf{u})}) \} + \mu_{q(1/a_\varepsilon)} \\ \mu_{q(1/\sigma_\varepsilon^2)} & \leftarrow \frac{1}{2} \left(o \sum_{i=1}^m n_i + 1 \right) / B_{q(\sigma_\varepsilon^2)}; \quad B_{q(a_\varepsilon)} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1/B_{q(a_\varepsilon)} \\ B_{q(\sigma_{\text{gbl}}^2)} & \leftarrow \frac{1}{2} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{gbl}})}) \} + \mu_{q(1/a_{\text{gbl}})} \\ \mu_{q(1/\sigma_{\text{gbl}}^2)} & \leftarrow \frac{1}{2} (L_{\text{gbl}} + 1) / B_{q(\sigma_{\text{gbl}}^2)} \\ B_{q(a_{\text{gbl}})} & \leftarrow \mu_{q(1/\sigma_{\text{gbl}}^2)} + (A_{\text{gbl}})^{-2}; \quad \mu_{q(1/a_{\text{gbl}})} \leftarrow 1/B_{q(a_{\text{gbl}})} \\ B_{q(\sigma_{\text{grpout}}^2)} & \leftarrow \frac{1}{2} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grpout}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{grpout}})}) \} + \mu_{q(1/a_{\text{grpout}})} \\ \mu_{q(1/\sigma_{\text{grpout}}^2)} & \leftarrow \frac{1}{2} (m L_{\text{grpout}} + 1) / B_{q(\sigma_{\text{grpout}}^2)} \\ B_{q(a_{\text{grpout}})} & \leftarrow \mu_{q(1/\sigma_{\text{grpout}}^2)} + A_{\text{grpout}}^{-2}; \quad \mu_{q(1/a_{\text{grpout}})} \leftarrow 1/B_{q(a_{\text{grpout}})} \\ \mathbf{B}_{q(\Sigma_{\text{grpout}})} & \leftarrow \sum_{i=1}^m \left\{ \boldsymbol{\mu}_{q(\delta_i)} \boldsymbol{\mu}_{q(\delta_i)}^\top + \Sigma_{q(\delta_i)} \right\} \\ & \quad + 2\nu_{\text{grpout}} \text{diag} \left(\mu_{q(1/a_{\Sigma_{\text{grpout},1}})}, \mu_{q(1/a_{\Sigma_{\text{grpout},2}})} \right) \\ \mathbf{M}_{q(\Sigma_{\text{grpout}}^{-1})} & \leftarrow (\nu_{\text{grpout}} + m + 1) \mathbf{B}_{q(\Sigma_{\text{grpout}})}^{-1} \\ \text{For } j = 1, 2: \quad B_{q(a_{\Sigma_{\text{grpout},j}})} & \leftarrow \nu_{\text{grpout}} \left(\mathbf{M}_{q(\Sigma_{\text{grpout}}^{-1})} \right)_{jj} + 1/A_{\Sigma_{\text{grpout},j}}^2 \\ \mu_{q(1/a_{\Sigma_{\text{grpout},j}})} & \leftarrow \left(\frac{\nu_{\text{grpout}}}{2} + 1 \right) / B_{q(a_{\Sigma_{\text{grpout},j}})} \\ B_{q(\sigma_{\text{grpinn}}^2)} & \leftarrow \frac{1}{2} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grpinn}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{grpinn}})}) \} + \mu_{q(1/a_{\text{grpinn}})} \\ \mu_{q(1/\sigma_{\text{grpinn}}^2)} & \leftarrow \frac{1}{2} (om L_{\text{grpinn}} + 1) / B_{q(\sigma_{\text{grpinn}}^2)} \\ B_{q(a_{\text{grpinn}})} & \leftarrow \mu_{q(1/\sigma_{\text{grpinn}}^2)} + A_{\text{grpinn}}^{-2}; \quad \mu_{q(1/a_{\text{grpinn}})} \leftarrow 1/B_{q(a_{\text{grpinn}})} \end{aligned}$$

$$\begin{aligned}
 \mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{grpinn}})} &\leftarrow \sum_{i=1}^m \sum_{j=1}^{n_i} \left\{ \boldsymbol{\mu}_{q(\gamma_{ij})} \boldsymbol{\mu}_{q(\gamma_{ij})}^\top + \boldsymbol{\Sigma}_{q(\gamma_{ij})} \right\} \\
 &\quad + 2\nu_{\text{grpinn}} \text{diag} \left(\mu_{q(1/a_{\Sigma_{\text{grpinn},1}})}, \mu_{q(1/a_{\Sigma_{\text{grpinn},2}})} \right) \\
 \mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpinn}}^{-1})} &\leftarrow \left(\nu_{\text{grpinn}} + \sum_{i=1}^m n_i + 1 \right) \mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{grpinn}})}^{-1} \\
 \text{For } j = 1, 2: \quad \mathbf{B}_{q(a_{\Sigma_{\text{grpinn},j}})} &\leftarrow \nu_{\text{grpinn}} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpinn}}^{-1})} \right)_{jj} + 1/A_{\Sigma_{\text{grpinn},j}}^2 \\
 \mu_{q(1/a_{\Sigma_{\text{grpinn},j}})} &\leftarrow \left(\frac{\nu_{\text{grpinn}}}{2} + 1 \right) / B_{q(a_{\Sigma_{\text{grpinn},j}})}
 \end{aligned}$$

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

Algorithm 7: *Naïve MFVB algorithm for the estimation of the optimal parameters in (5.21).*

is shown that $\log \{ \underline{p}(\mathbf{y}; q) \}$ has the following expression:

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) = & \\
 & \frac{1}{2}(\nu_{\text{grpout}} + q_{\text{grpout}} - 1) \log(2\nu_{\text{grpout}}) + \frac{1}{2}(\nu_{\text{grpinn}} + q_{\text{grpinn}} - 1) \log(2\nu_{\text{grpinn}}) \\
 & - \frac{1}{2} o \sum_{i=1}^m n_i \log(2\pi) - 5 \log(\pi) - \frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} \\
 & + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| + \log \Gamma \left\{ \frac{1}{2} (L_{\text{gbl}} + 1) \right\} - \log (C_{q_{\text{grpout}}, \nu_{\text{grpout}} + q_{\text{grpout}} - 1}) \\
 & + \frac{1}{2} \left\{ p + 2L_{\text{grpout}} + m(q_{\text{grpout}} + L_{\text{grpout}}) + \sum_{i=1}^m n_i (q_{\text{grpinn}} + L_{\text{grpinn}}) \right\} - \log (B_{q(a_{\text{gbl}})}) \\
 & + \log (C_{q_{\text{grpout}}, \nu_{\text{grpout}} + m + q_{\text{grpout}} - 1}) - \log (C_{q_{\text{grpinn}}, \nu_{\text{grpinn}} + q_{\text{grpinn}} - 1}) \\
 & + \log (C_{q_{\text{grpinn}}, \nu_{\text{grpinn}} + \sum_{i=1}^m n_i + q_{\text{grpinn}} - 1}) + \log \Gamma \left\{ \frac{1}{2} \left(o \sum_{i=1}^m n_i + 1 \right) \right\} - \log (A_\varepsilon) \\
 & - \frac{1}{2} (\nu_{\text{grpout}} + q_{\text{grpout}} + m - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{grpout}})}| - \log (B_{q(a_\varepsilon)}) - \log (A_{\text{gbl}}) \\
 & - \frac{1}{2} \left(\nu_{\text{grpinn}} + q_{\text{grpinn}} + \sum_{i=1}^m n_i - 1 \right) \log |\mathbf{B}_{q(\boldsymbol{\Sigma}_{\text{grpinn}})}| - \frac{1}{2} \left(o \sum_{i=1}^m n_i + 1 \right) \log (B_{q(\sigma_\varepsilon^2)}) \\
 & + \log \Gamma \left\{ \frac{1}{2} \left(L_{\text{grpinn}} \sum_{i=1}^m n_i + 1 \right) \right\} - \frac{1}{2} (L_{\text{gbl}} + 1) \log (B_{q(\sigma_{\text{gbl}}^2)}) - \log (A_{\text{grpinn}}) \\
 & - \frac{1}{2} (mL_{\text{grpout}} + 1) \log (B_{q(\sigma_{\text{grpout}}^2)}) - \frac{1}{2} \left(L_{\text{grpinn}} \sum_{i=1}^m n_i + 1 \right) \log (B_{q(\sigma_{\text{grpinn}}^2)}) \\
 & + \mu_{q(1/\sigma_\varepsilon^2)} \mu_{q(1/a_\varepsilon)} + \mu_{q(1/\sigma_{\text{gbl}}^2)} \mu_{q(1/a_{\text{gbl}})} + \log \Gamma \left\{ \frac{1}{2} (mL_{\text{grpout}} + 1) \right\} \\
 & - \log (B_{q(a_{\text{grpout}})}) + \mu_{q(1/\sigma_{\text{grpout}}^2)} \mu_{q(1/a_{\text{grpout}})} - \sum_{j=1}^{q_{\text{grpinn}}} \log (A_{\Sigma_{\text{grpinn},j}}) \\
 & - \log (A_{\text{grpout}}) - \sum_{j=1}^{q_{\text{grpout}}} \log (A_{\Sigma_{\text{grpout},j}}) + q_{\text{grpout}} \log \Gamma \left\{ \frac{1}{2} (\nu_{\text{grpout}} + q_{\text{grpout}}) \right\} \\
 & - \frac{1}{2} (\nu_{\text{grpout}} + q_{\text{grpout}} - 1) \sum_{j=1}^{q_{\text{grpout}}} \log (B_{q(a_{\Sigma_{\text{grpout},j}})}) + q_{\text{grpinn}} \log \Gamma \left\{ \frac{1}{2} (\nu_{\text{grpinn}} + q_{\text{grpinn}}) \right\}
 \end{aligned} \tag{5.22}$$

$$\begin{aligned}
& + \sum_{j=1}^{q_{\text{grpout}}} \nu_{\text{grpout}} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpout}}^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma_{S \cdot j}})} - \log \left(B_{q(a_{\text{grpinn}})} \right) \\
& - \frac{1}{2} (\nu_{\text{grpinn}} + q_{\text{grpinn}} - 1) \sum_{j=1}^{q_{\text{grpinn}}} \log \left(B_{q(a_{\Sigma_{\text{grpinn} \cdot j}})} \right) + \mu_{q(1/\sigma_{\text{grpinn}}^2)} \mu_{q(1/a_{\text{grpinn}})} \\
& + \sum_{j=1}^{q_{\text{grpinn}}} \nu_{\text{grpinn}} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}_{\text{grpinn}}^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma_{\text{grpinn} \cdot j}})}.
\end{aligned}$$

where $p = 6$ is the number of columns in the \mathbf{X} matrix, $q_{\text{grpout}} = 2$ is the dimension of $\boldsymbol{\Sigma}_{\text{grpout}}$ and $q_{\text{grpinn}} = 2$ is the dimension of $\boldsymbol{\Sigma}_{\text{grpinn}}$.

5.3.3 Streamlining mean field variational Bayes for the three-level Gaussian response model

Here we extend the streamlined MFVB algorithm for the two-level Gaussian response model as shown in Section 5.2.4 to cater to our three-level structured rodent tumor data. Once again, in keeping with the notation in Lee & Wand (2015) we partition the following vectors and matrices:

$$\mathbf{Z}^{\text{R}} \equiv \text{blockdiag} \left(\mathbf{Z}_i^{\text{R}} \right)_{1 \leq i \leq m},$$

where

$$\mathbf{Z}_i^{\text{R}} \equiv \left[\mathbf{X}_i \mid \mathbf{Z}_i^{\text{grpout}} \mid \text{blockdiag} \left(\mathbf{X}_{ij} \mathbf{Z}_{ij}^{\text{grpinn}} \right)_{1 \leq i \leq m} \right].$$

The global and random effects are partitioned according to:

$$\mathbf{u} \equiv \begin{bmatrix} \mathbf{u}^{\text{gbl}} \\ \mathbf{u}^{\text{R}} \end{bmatrix}, \quad \mathbf{Z} \equiv [\mathbf{Z}^{\text{gbl}} \quad \mathbf{Z}^{\text{R}}],$$

where

$$\mathbf{u}^{\text{R}} \equiv \begin{bmatrix} \mathbf{u}_1^{\text{R}} \\ \vdots \\ \mathbf{u}_m^{\text{R}} \end{bmatrix}, \quad \mathbf{u}_i^{\text{R}} = \begin{bmatrix} \delta_i \\ \mathbf{u}_i^{\text{grpout}} \\ \gamma_{i1} \\ \mathbf{u}_{i1}^{\text{grpinn}} \\ \vdots \\ \gamma_{in_i} \\ \mathbf{u}_{in_i}^{\text{grpinn}} \end{bmatrix}.$$

With these partitions in mind, a streamlined version of Algorithm 7 is presented in Algorithm 8. Successive values of the marginal log-likelihood lower bound are used to diagnose convergence of Algorithm 8. These values are monitored using the expression in (5.22) where

$$\log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|$$

Set up initial values:

$$\begin{aligned} & \mu_q(1/\sigma_\varepsilon^2), \mu_q(1/\sigma_{\text{gbl}}^2), \mu_q(1/\sigma_{\text{grpout}}^2), \mu_q(1/\sigma_{\text{grpinn}}^2), \mu_q(1/a_\varepsilon), \mu_q(1/a_{\text{gbl}}), \mu_q(1/a_{\text{grpout}}), \\ & \mu_q(1/a_{\text{grpinn}}) > 0, \mu_q(1/a_{\Sigma_{\text{grpout}}:j}), \mu_q(1/a_{\Sigma_{\text{grpinn}}:j}) > 0, 1 \leq j \leq 2, \mathbf{M}_{q(\Sigma_{\text{grpout}}^{-1})}, \\ & \mathbf{M}_{q(\Sigma_{\text{grpinn}}^{-1})} \text{ positive definite.} \end{aligned}$$

Cycle through:

$$\mathbf{S} \leftarrow \mathbf{0}; \quad \mathbf{s} \leftarrow \mathbf{0}$$

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

$$\begin{aligned} \mathbf{G}_i & \leftarrow \mu_q(1/\sigma_\varepsilon^2) (\mathbf{C}_i^{\text{gbl}})^\top \mathbf{Z}_i^{\text{R}} \\ \mathbf{H}_i & \leftarrow \left[\mu_q(1/\sigma_\varepsilon^2) (\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} + \text{blockdiag} \left\{ M_{q(\Sigma_{\text{grpout}}^{-1})}, \mu_q(1/\sigma_{\text{grpout}}^2) \mathbf{I}_{L_{\text{grpout}}}, \right. \right. \\ & \left. \left. \mathbf{I}_{n_i} \otimes \left(M_{q(\Sigma_{\text{grpinn}}^{-1})}, \mu_q(1/\sigma_{\text{grpinn}}^2) \mathbf{I}_{L_{\text{grpinn}}} \right) \right\} \right]^{-1} \\ \mathbf{S} & \leftarrow \mathbf{S} + \mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top; \quad \mathbf{s} \leftarrow \mathbf{s} + \mathbf{G}_i \mathbf{H}_i (\mathbf{Z}_i^{\text{R}})^\top \mathbf{y}_i \end{aligned}$$

$$\Sigma_{q(\beta, \mathbf{u}^{\text{G}})} \leftarrow \left\{ \mu_q(1/\sigma_\varepsilon^2) (\mathbf{C}^{\text{gbl}})^\top \mathbf{C}^{\text{gbl}} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_6 & \mathbf{0} \\ \mathbf{0} & \mu_q(1/\sigma_{\text{gbl}}^2) \mathbf{I}_{L_{\text{gbl}}} \end{bmatrix} - \mathbf{S} \right\}^{-1}$$

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})} \leftarrow \mu_q(1/\sigma_\varepsilon^2) \Sigma_{q(\beta, \mathbf{u}^{\text{G}})} \{ (\mathbf{C}^{\text{gbl}})^\top \mathbf{y} - \mathbf{s} \}$$

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

$$\begin{aligned} \Sigma_{q(\mathbf{u}_i^{\text{R}})} & \leftarrow \mathbf{H}_i + \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^{\text{G}})} \mathbf{G}_i \mathbf{H}_i \\ \boldsymbol{\mu}_{q(\mathbf{u}_i^{\text{R}})} & \leftarrow \mathbf{H}_i \{ \mu_q(1/\sigma_\varepsilon^2) (\mathbf{Z}_i^{\text{R}})^\top \mathbf{y}_i - \mathbf{G}_i^\top \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})} \} \\ B_{q(\sigma_\varepsilon^2)} & \leftarrow \mu_q(1/a_\varepsilon) + \frac{1}{2} \left[\left\| \mathbf{y} - \mathbf{C}^{\text{gbl}} \boldsymbol{\mu}_{q(\beta, \mathbf{u}^{\text{G}})} - \begin{bmatrix} \mathbf{Z}_1^{\text{R}} \boldsymbol{\mu}_{q(\mathbf{u}_1^{\text{R}})} \\ \vdots \\ \mathbf{Z}_m^{\text{R}} \boldsymbol{\mu}_{q(\mathbf{u}_m^{\text{R}})} \end{bmatrix} \right\|^2 \right. \\ & \quad \left. + \text{tr} \{ (\mathbf{C}^{\text{gbl}})^\top \mathbf{C}^{\text{gbl}} \Sigma_{q(\beta, \mathbf{u}^{\text{G}})} \} + \sum_{i=1}^m \text{tr} \{ (\mathbf{Z}_i^{\text{R}})^\top \mathbf{Z}_i^{\text{R}} \Sigma_{q(\mathbf{u}_i^{\text{R}})} \} \right. \\ & \quad \left. - 2 \mu_q^{-1}(1/\sigma_\varepsilon^2) \sum_{i=1}^m \text{tr} (\mathbf{G}_i \mathbf{H}_i \mathbf{G}_i^\top \Sigma_{q(\beta, \mathbf{u}^{\text{G}})}) \right] \end{aligned}$$

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

$$B_{q(\sigma_{\text{gbl}}^2)} \leftarrow \frac{1}{2} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{gbl}})}) \} + \mu_q(1/a_{\text{gbl}})$$

$$\mu_q(1/\sigma_{\text{gbl}}^2) \leftarrow \frac{1}{2} (L_{\text{gbl}} + 1) / B_{q(\sigma_{\text{gbl}}^2)}$$

$$B_{q(a_{\text{gbl}})} \leftarrow \mu_q(1/\sigma_{\text{gbl}}^2) + (A_{\text{gbl}})^{-2}; \quad \mu_q(1/a_{\text{gbl}}) \leftarrow 1 / B_{q(a_{\text{gbl}})}$$

$$B_{q(\sigma_{\text{grpout}}^2)} \leftarrow \frac{1}{2} \{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{grpout}})}\|^2 + \text{tr}(\Sigma_{q(\mathbf{u}^{\text{grpout}})}) \} + \mu_q(1/a_{\text{grpout}})$$

$$\begin{aligned}\mu_q(1/\sigma_{\text{grpout}}^2) &\leftarrow \frac{1}{2} (m L_{\text{grpout}} + 1) / B_q(\sigma_{\text{grpout}}^2) \\ B_q(a_{\text{grpout}}) &\leftarrow \mu_q(1/\sigma_{\text{grpout}}^2) + A_{\text{grpout}}^{-2}, \quad \mu_q(1/a_{\text{grpout}}) \leftarrow 1/B_q(a_{\text{grpout}}) \\ B_q(\Sigma_{\text{grpout}}) &\leftarrow \sum_{i=1}^m \left\{ \boldsymbol{\mu}_q(\boldsymbol{\delta}_i) \boldsymbol{\mu}_q^\top(\boldsymbol{\delta}_i) + \Sigma_q(\boldsymbol{\delta}_i) \right\} \\ &\quad + 2\nu_{\text{grpout}} \text{diag} \left(\mu_q(1/a_{\Sigma_{\text{grpout},1}}), \mu_q(1/a_{\Sigma_{\text{grpout},2}}) \right) \\ \mathbf{M}_q(\Sigma_{\text{grpout}}^{-1}) &\leftarrow (\nu_{\text{grpout}} + m + 1) \mathbf{B}_q^{-1}(\Sigma_{\text{grpout}})\end{aligned}$$

For $j = 1, 2$:

$$\begin{aligned}B_q(a_{\Sigma_{\text{grpout},j}}) &\leftarrow \nu_{\text{grpout}} \left(\mathbf{M}_q(\Sigma_{\text{grpout}}^{-1}) \right)_{jj} + 1/A_{\Sigma_{\text{grpout},j}}^2 \\ \mu_q(1/a_{\Sigma_{\text{grpout},j}}) &\leftarrow \left(\frac{\nu_{\text{grpout}}}{2} + 1 \right) / B_q(a_{\Sigma_{\text{grpout},j}}) \\ B_q(\sigma_{\text{grpinn}}^2) &\leftarrow \frac{1}{2} \left\{ \|\boldsymbol{\mu}_q(\mathbf{u}_{\text{grpinn}})\|^2 + \text{tr}(\Sigma_q(\mathbf{u}_{\text{grpinn}})) \right\} + \mu_q(1/a_{\text{grpinn}}) \\ \mu_q(1/\sigma_{\text{grpinn}}^2) &\leftarrow \frac{1}{2} (om L_{\text{grpinn}} + 1) / B_q(\sigma_{\text{grpinn}}^2) \\ B_q(a_{\text{grpinn}}) &\leftarrow \mu_q(1/\sigma_{\text{grpinn}}^2) + A_{\text{grpinn}}^{-2}, \quad \mu_q(1/a_{\text{grpinn}}) \leftarrow 1/B_q(a_{\text{grpinn}}) \\ B_q(\Sigma_{\text{grpinn}}) &\leftarrow \sum_{i=1}^m \sum_{j=1}^{n_i} \left\{ \boldsymbol{\mu}_q(\boldsymbol{\gamma}_{ij}) \boldsymbol{\mu}_q^\top(\boldsymbol{\gamma}_{ij}) + \Sigma_q(\boldsymbol{\gamma}_{ij}) \right\} \\ &\quad + 2\nu_{\text{grpinn}} \text{diag} \left(\mu_q(1/a_{\Sigma_{\text{grpinn},1}}), \mu_q(1/a_{\Sigma_{\text{grpinn},2}}) \right) \\ \mathbf{M}_q(\Sigma_{\text{grpinn}}^{-1}) &\leftarrow \left(\nu_{\text{grpinn}} + \sum_{i=1}^m n_i + 1 \right) \mathbf{B}_q^{-1}(\Sigma_{\text{grpinn}})\end{aligned}$$

For $j = 1, 2$:

$$\begin{aligned}B_q(a_{\Sigma_{\text{grpinn},j}}) &\leftarrow \nu_{\text{grpinn}} \left(\mathbf{M}_q(\Sigma_{\text{grpinn}}^{-1}) \right)_{jj} + 1/A_{\Sigma_{\text{grpinn},j}}^2 \\ \mu_q(1/a_{\Sigma_{\text{grpinn},j}}) &\leftarrow \left(\frac{\nu_{\text{grpinn}}}{2} + 1 \right) / B_q(a_{\Sigma_{\text{grpinn},j}})\end{aligned}$$

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

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

$$\begin{aligned}\boldsymbol{\Lambda}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}, \mathbf{u}_i^{\text{R}}) &\equiv E_q \left[\left\{ \left[\begin{array}{c} \boldsymbol{\beta} \\ \mathbf{u}^{\text{gbl}} \end{array} \right] - \boldsymbol{\mu}_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}) \right\} \left\{ \mathbf{u}_i^{\text{R}} - \boldsymbol{\mu}_q(\mathbf{u}_i^{\text{R}}) \right\}^\top \right] \\ &\leftarrow -\Sigma_q(\boldsymbol{\beta}, \mathbf{u}^{\text{gbl}}) \mathbf{G}_i \mathbf{H}_i\end{aligned}$$

Algorithm 8: Streamlined MFVB algorithm for the three-level Gaussian response model given in (5.19).

is replaced by

$$-\sum_{i=1}^m \log \left| \mu_{q(1/\sigma_{\varepsilon}^2)} (\mathbf{Z}_i^R)^\top \mathbf{Z}_i^R + \text{blockdiag} \left\{ M_{q(\Sigma^{-1})}, \mu_{q(1/\sigma_{\text{grpout}}^2)} \mathbf{I}_{L_{\text{grpout}}}, \right. \right. \\ \left. \left. \mathbf{I}_{n_i} \otimes \left(M_{q(\Sigma_{\text{grpinn}}^{-1})}, \mu_{q(1/\sigma_{\text{grpinn}}^2)} \mathbf{I}_{L_{\text{grpinn}}} \right) \right\} \right| - \log |\Sigma_{q(\beta, \mathbf{u}^{\text{gbl}})}|.$$

Once again, Algorithm 8 should not be considered a fully streamlined version of Algorithm 7 because of the structure of \mathbf{H}_i . In addition to incorporating the variance structure of the spline components for each subject's curve, there is also the variance structure associated with the slice specific curves. This is represented in \mathbf{H}_i as $\mathbf{I}_{\sum_{i=1}^m n_i} \otimes \left(M_{q(\Sigma_{\text{grpinn}}^{-1})}, \mu_{q(1/\sigma_{\text{grpinn}}^2)} \mathbf{I}_{L_{\text{grpinn}}} \right)$. This block diagonal structure, comprising storage of a large amount of zeros, is hindering the potential time savings that would be possible given further research into this area.

5.3.3.1 Results

The three-level model given in (5.19) was fit using the streamlined MFVB algorithm as shown in Algorithm 8. For fitting this data we set $L_{\text{gbl}} = 17$, $L_{\text{grpout}} = 12$ and $L_{\text{grpinn}} = 9$ and set all hyperparameters to 10^5 . The `rstan` package was used for performing MCMC with a burn-in of 1000 and 1000 additional iterations. The MFVB iterations were terminated when the increase in the corresponding $\log p(\mathbf{b}_{\text{sc}}; q)$ fell below 10^{-8} . The MCMC fit took approximately 18 hours whilst the streamlined MFVB fit took just under 2 hours to run. Figure 5.8 illustrates the accuracy of the streamlined MFVB approach against the MCMC benchmark. The parameters being monitored are the quartiles of the curve estimates of the 1st rodent for three of its slice curves. These are represented as $f(Q_i) + g_1(Q_i) + h_{11}(Q_i)$, $f(Q_i) + g_1(Q_i) + h_{12}(Q_i)$ and $f(Q_i) + g_1(Q_i) + h_{13}(Q_i)$ for $i = 1, 2, 3$. The accuracy of the curve estimates for these parameters are extremely high, especially considering the complexity of the model.

Figures 5.9 and 5.10 provide illustration of the inference provided by the streamlined MFVB approach and MCMC for the three-level Gaussian response model. Each of these two figures represents a rodent with the interior panels corresponding to the slices or probe locations in which each tumor was scanned. The streamlined MFVB fits are almost indistinguishable from the MCMC fits for each slice.

5.3. THREE-LEVEL GAUSSIAN RESPONSE MODEL

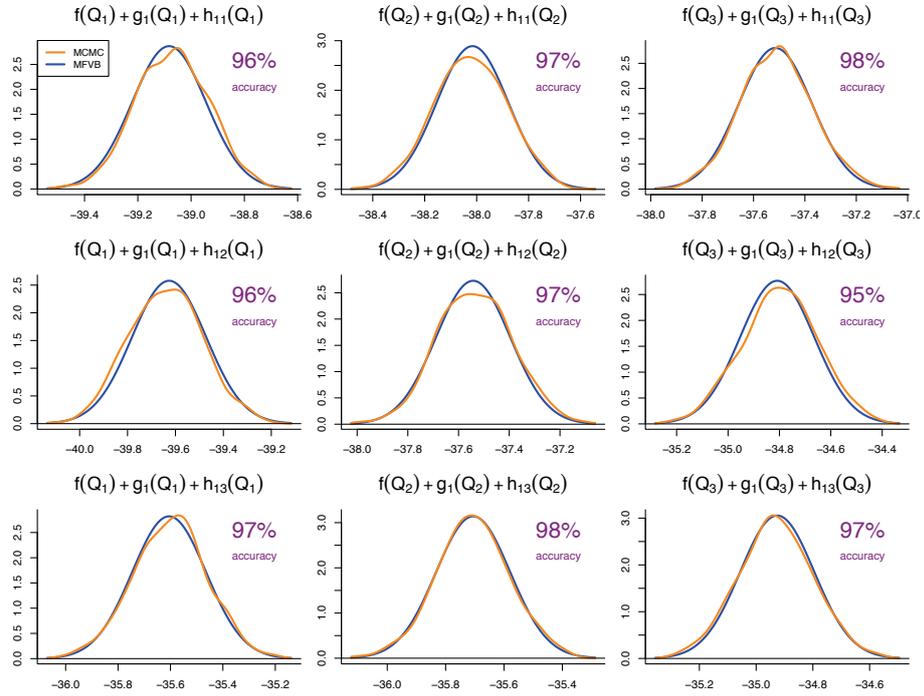


Figure 5.8: Approximate posterior density functions obtained via MCMC and streamlined MFVB for the quartiles of three specific slice curves from the rodent tumor data-set. MFVB accuracy scores are displayed.

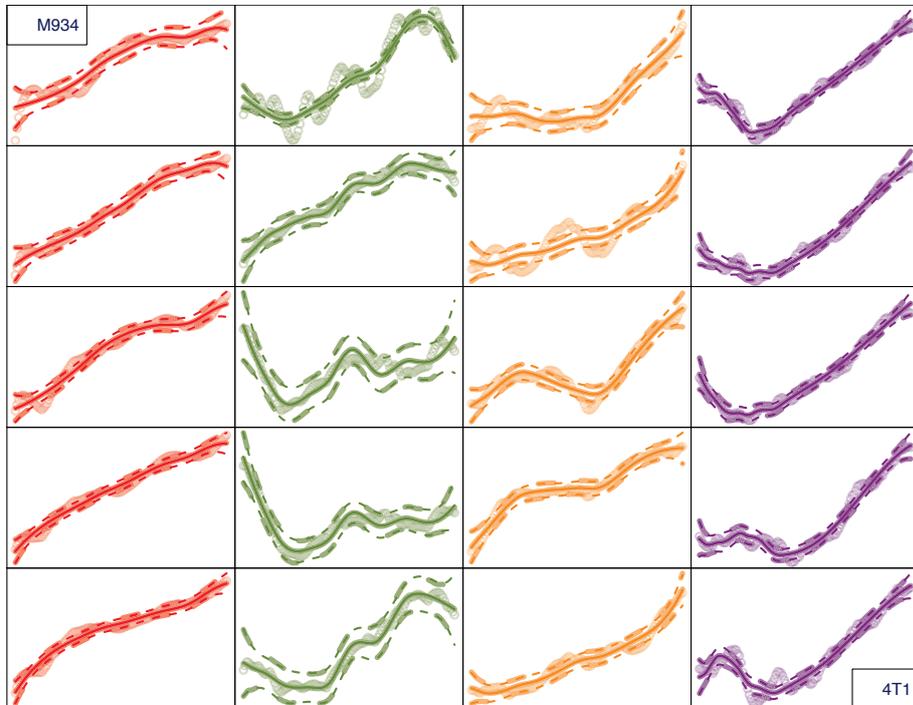


Figure 5.9: Slice-specific estimates and pointwise 95 % credible sets for a rodent given in Figure 5.6. The different colours correspond to the different transducers. The thicker curve estimates correspond to MCMC and the thinner ones correspond to streamlined MFVB fitting. The data points for each slice are also plotted.

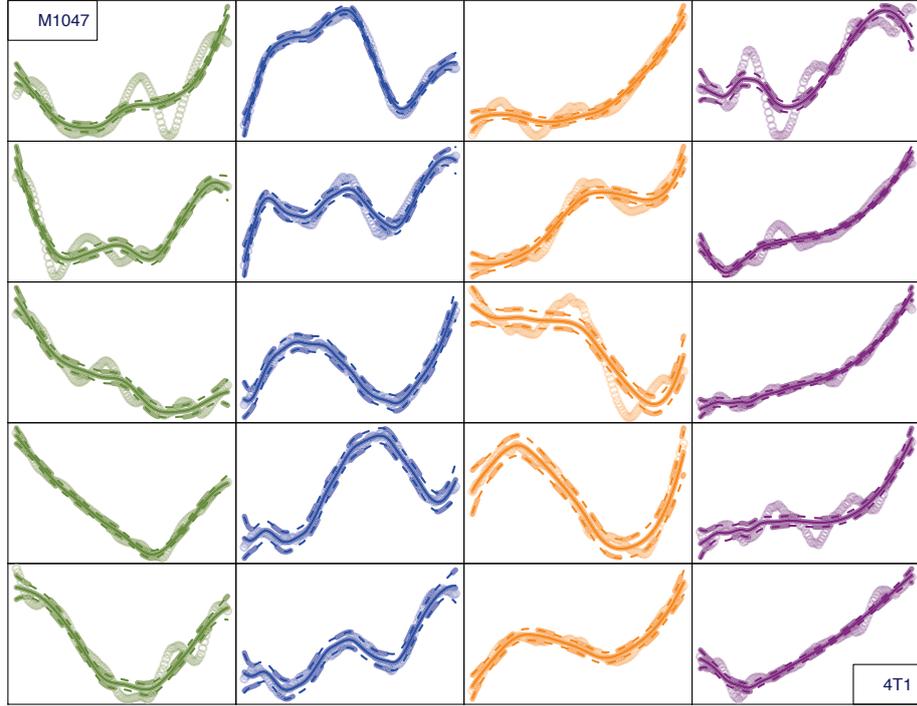


Figure 5.10: *Slice-specific estimates and pointwise 95 % credible sets for a rodent given in Figure 5.6. The different colours correspond to the different transducers. The thicker curve estimates correspond to MCMC and the thinner ones correspond to streamlined MFVB fitting. The data points for each slice are also plotted.*

5.3.4 Simulation study

An extensive simulation study was carried out to assess the speed of Algorithm 8 against that of Algorithm 7. We generated 25 data-sets according to

$$y_{ijk} = f(x_{ijk}) + g_i(x_{ijk}) + h_{ij}(x_{ijk}) + \varepsilon_{ijk}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n_i, \quad 1 \leq k \leq o_{ij},$$

where the $x_{ijk} \stackrel{\text{ind.}}{\sim} \text{Uniform}(0, 1)$, $\varepsilon_{ijk} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_\varepsilon^2)$, $u_\ell^{\text{gbl}} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\text{gbl}}^2)$, $\delta_i \stackrel{\text{ind.}}{\sim} \text{N}(\mathbf{0}, \Sigma_{\text{grpout}})$, $u_{i\ell}^{\text{grpout}} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\text{grpout}}^2)$, $\gamma_{ij} \stackrel{\text{ind.}}{\sim} \text{N}(\mathbf{0}, \Sigma_{\text{grpinn}})$, and $u_{ij\ell}^{\text{grpinn}} \stackrel{\text{ind.}}{\sim} \text{N}(0, \sigma_{\text{grpinn}}^2)$. The true parameter values were specified as

$$\beta_0 = 0.2, \quad \beta_1 = 1.8, \quad \sigma_\varepsilon^2 = 0.02, \quad \sigma_{\text{gbl}}^2 = 1.5, \quad \Sigma_{\text{grpout}} = \begin{bmatrix} 0.4 & 0.1 \\ 0.1 & 0.4 \end{bmatrix}, \quad \sigma_{\text{grpout}}^2 = 0.29,$$

$$\Sigma_{\text{grpinn}} = \begin{bmatrix} 0.4 & 0.15 \\ 0.15 & 0.4 \end{bmatrix} \quad \text{and} \quad \sigma_{\text{grpinn}}^2 = 0.17.$$

The number of subjects between simulation studies was varied, where $m \in \{5, 15, 25, 35\}$, the within subject curves remained constant at $n_i = 5$, $1 \leq i \leq m$ and the within curve

5.4. DISCUSSION

sample size remained constant at $o_{ij} = 10$, $1 \leq i \leq m$, $1 \leq j \leq n_i$. The computations were performed on a laptop computer (Mac OS X; 2.8 GHz processor, 16 GBytes of random access memory). The computation times for each approach are summarised in Table 5.4. The time savings here will not be as impressive as the time savings given in Table 5.1,

m	Naïve	Streamlined	Ratio
5	28.49 (3.98)	4.34 (0.15)	6.57
15	617.65 (11.37)	44.79 (1.77)	13.79
25	2638.30 (20.81)	177.42 (0.19)	14.87
35	7051.99 (22.47)	455.89 (2.06)	15.47

Table 5.4: *Average (standard deviation) run time in seconds for naïve and streamlined MFVB fitting of the three-level Gaussian response model.*

due to the storage issues contained in \mathbf{H}_i . However, we still notice that the streamlined approach for the three-level Gaussian response model is significantly faster than the Naïve approach.

5.4 Discussion

The original aim of this chapter was to develop a fast deterministic approach to MCMC to shorten estimation time for the two-level and three-level Gaussian response model. The rodent tumor dataset that was introduced in Section 5.3 led to the consideration of a possible streamlined approach to MFVB. Meanwhile, research was being conducted by Lee & Wand (2015) that led to the fully streamlined MFVB algorithm for a similar two-level Gaussian response model. This led to the incorporation of a variant of this streamlined algorithm for this chapter. As is evident from the real datasets and simulation studies for both the two-level and three-level Gaussian response models considered here, MFVB in general has achieved faster inference, but a streamlined version of MFVB in particular has achieved even faster inference with little loss in accuracy. The models considered in this chapter are of reasonably high complexity and further research into a fully streamlined algorithm for both models would indeed be beneficial.

5.A Derivation of Algorithm 5

Expressions for $\mu_{q(\beta, \mathbf{u})}$ and $\Sigma_{q(\beta, \mathbf{u})}$

Firstly, we define

$$\mathbf{G} \equiv \begin{bmatrix} \sigma_{\text{gbl}}^2 \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\boldsymbol{\Sigma}, \sigma_{\text{grp}}^2 \mathbf{I}_{L_{\text{grp}}}) \end{bmatrix}.$$

Next, we note that

$$\begin{aligned} p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_{\varepsilon}^2) p(\mathbf{u} | \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2) p(\boldsymbol{\beta}) \\ &\propto \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u})^\top (\sigma_{\varepsilon}^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}) \right\} \times \exp \left\{ -\frac{1}{2} \mathbf{u}^\top \mathbf{G}^{-1} \mathbf{u} \right\} \\ &\quad \times \exp \left(-\frac{1}{2\sigma_{\beta}^2} \|\boldsymbol{\beta}\|^2 \right). \end{aligned}$$

We next combine the design matrices \mathbf{X} and \mathbf{Z} to form $\mathbf{C} = [\mathbf{X} | \mathbf{Z}]$. Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) &\propto -\frac{1}{2} \left(\mathbf{y} - \mathbf{C} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right)^\top (\sigma_{\varepsilon}^{-2} \mathbf{I}_n) \left(\mathbf{y} - \mathbf{C} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right) - \frac{1}{2} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \begin{bmatrix} \sigma_{\beta}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}^{-1} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \\ &= -\frac{1}{2} \left\{ \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \left(\mathbf{C}^\top (\sigma_{\varepsilon}^{-2} \mathbf{I}_n) \mathbf{C} + \begin{bmatrix} \sigma_{\beta}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}^{-1} \end{bmatrix} \right) \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} - 2 \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \mathbf{C}^\top (\sigma_{\varepsilon}^{-2} \mathbf{I}_n) \mathbf{y} \right\} \\ &\quad + \text{const} \end{aligned}$$

where ‘const’ denotes terms not depending on the parameter vector $(\boldsymbol{\beta}, \mathbf{u})$. Then, taking expectations with respect to all parameters except $(\boldsymbol{\beta}, \mathbf{u})$:

$$\begin{aligned} \log q^*(\boldsymbol{\beta}, \mathbf{u}) &= E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \text{rest}) \} + \text{const} \\ &= -\frac{1}{2} \left\{ \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \left(\mathbf{C}^\top (\mu_{q(1/\sigma_{\varepsilon}^2)} \mathbf{I}_n) \mathbf{C} + \begin{bmatrix} \sigma_{\beta}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{M}_{q(\mathbf{G}^{-1})} \end{bmatrix} \right) \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right. \\ &\quad \left. - 2 \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}^\top \mathbf{C}^\top (\mu_{q(1/\sigma_{\varepsilon}^2)} \mathbf{I}_n) \mathbf{y} \right\} + \text{const} \end{aligned}$$

where

$$\mathbf{M}_{q(\mathbf{G}^{-1})} \equiv \begin{bmatrix} \mu_{q(1/\sigma_{\text{gbl}}^2)} \mathbf{I}_{L_{\text{gbl}}} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq i \leq m} \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})}, \mu_{q(1/\sigma_{\text{grp}}^2)} \mathbf{I}_{L_{\text{grp}}} \right) \end{bmatrix},$$

and

$$\mathbf{M}_q(\mathbf{X}) \equiv E_q(\mathbf{X}).$$

Completing the square from above gives

$$q^* \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right) \sim \text{N}(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$$

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \mu_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top \mathbf{y}, \text{ and} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} &= \left(\mu_{q(1/\sigma_\varepsilon^2)} \mathbf{C}^\top \mathbf{C} + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{M}_{q(\mathbf{G}^{-1})} \end{bmatrix} \right)^{-1} \end{aligned}$$

where $p = 2$ is the column length in \mathbf{X} .

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

The full conditional distribution for σ_ε^2 is given by

$$\begin{aligned} p(\sigma_\varepsilon^2 | \text{rest}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) p(\sigma_\varepsilon^2 | a_\varepsilon) \\ &= |\sigma_\varepsilon^2 \mathbf{I}_n|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u})^\top (\sigma_\varepsilon^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}) \right\} \\ &\quad (\sigma_\varepsilon^2)^{-3/2} \exp \left\{ -(1/a_\varepsilon) / \sigma_\varepsilon^2 \right\} + \text{const} \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\sigma_\varepsilon^2 | \text{rest}) &= -\frac{n}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2 - \frac{3}{2} \log(\sigma_\varepsilon^2) - (1/a_\varepsilon) / \sigma_\varepsilon^2 + \text{const} \\ &= \frac{1}{2} (n+3) \log(\sigma_\varepsilon^2) - \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2 + 1/a_\varepsilon \right) / \sigma_\varepsilon^2 + \text{const} \end{aligned}$$

Taking expectations:

$$\begin{aligned} \log q^*(\sigma_\varepsilon^2) &= E_q \left\{ \log p(\sigma_\varepsilon^2 | \text{rest}) \right\} + \text{const} \\ &= \frac{1}{2} (n+3) \log(\sigma_\varepsilon^2) - \left(\frac{1}{2} E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2 + \mu_{q(1/a_\varepsilon)} \right) / \sigma_\varepsilon^2 + \text{const} \end{aligned}$$

Using result 1.4.19 it follows that

$$\begin{aligned} \log q^*(\sigma_\varepsilon^2) &= \frac{1}{2}(n+3)\log(\sigma_\varepsilon^2) - \left(\frac{1}{2}\{\|\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})\} + \mu_{q(1/a_\varepsilon)}\right)/\sigma_\varepsilon^2 \\ &\quad + \text{const} \end{aligned}$$

This is of the form of an Inverse-Gamma distribution, therefore

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

where

$$B_{q(\sigma_\varepsilon^2)} = \frac{1}{2}\{\|\mathbf{y} - \mathbf{C}\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})\} + \mu_{q(1/a_\varepsilon)}$$

and using Result 1.4.3, we get

$$\mu_{q(1/\sigma_\varepsilon^2)} = \frac{1}{2}(n+1)/B_{q(\sigma_\varepsilon^2)}.$$

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

The full conditional distribution of a_ε is given by

$$\begin{aligned} p(a_\varepsilon | \text{rest}) &\propto p(\sigma_\varepsilon^2 | a_\varepsilon) p(a_\varepsilon) \\ &= (1/a_\varepsilon)^{1/2} \exp\{- (1/a_\varepsilon)/\sigma_\varepsilon^2\} \times (a_\varepsilon)^{-\frac{1}{2}-1} \exp\{- (1/A_\varepsilon^2)/a_\varepsilon\} + \text{const} \end{aligned}$$

Next, taking the logarithm gives

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

Taking expectations:

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

This is of the form of an Inverse-Gamma distribution, therefore

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

where

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

and using Result 1.4.3, we get

$$\mu_{q(1/a_\varepsilon)} = 1/B_{q(a_\varepsilon)}.$$

Expressions for $B_q(\sigma_{\text{gbl}}^2)$, and $\mu_{q(1/\sigma_{\text{gbl}}^2)}$

The full conditional distribution of σ_ε^2 is attained through

$$\begin{aligned} p(\sigma_{\text{gbl}}^2 | \text{rest}) &\propto p(\mathbf{u} | \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2) p(\sigma_{\text{gbl}}^2 | a_{\text{gbl}}) \\ &= |\mathbf{G}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \mathbf{u}^\top \mathbf{G}^{-1} \mathbf{u}\right\} (\sigma_{\text{gbl}}^2)^{-\frac{1}{2}-1} \exp\left\{(1/a_{\text{gbl}})/\sigma_{\text{gbl}}^2\right\} + \text{const} \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\sigma_{\text{gbl}}^2 | \text{rest}) &= -\frac{L_{\text{gbl}}}{2} \log(\sigma_{\text{gbl}}^2) - \frac{3}{2} \log(\sigma_{\text{gbl}}^2) - \frac{1}{\sigma_{\text{gbl}}^2} \|\mathbf{u}^{\text{gbl}}\|^2 - (1/a_{\text{gbl}})/\sigma_{\text{gbl}}^2 + \text{const} \\ &= -\frac{1}{2}(L_{\text{gbl}} + 3) \log(\sigma_{\text{gbl}}^2) - \left\{\frac{1}{2} \|\mathbf{u}^{\text{gbl}}\|^2 + 1/a_{\text{gbl}}\right\} / \sigma_{\text{gbl}}^2 + \text{const} \end{aligned}$$

Taking expectations with respect to all parameters but σ_{gbl}^2 , with the help of Result 1.4.17, we get

$$\begin{aligned} \log q^*(\sigma_{\text{gbl}}^2) &= E_q \left\{ \log p(\sigma_{\text{gbl}}^2 | \text{rest}) \right\} \\ &= -\frac{1}{2}(L_{\text{gbl}} + 3) \log(\sigma_{\text{gbl}}^2) - \left[\frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{gbl}})}) \right\} + \mu_{(1/a_{\text{gbl}})} \right] / \sigma_{\text{gbl}}^2 \\ &\quad + \text{const} \end{aligned}$$

Thus, we can represent the optimal q -density of σ_{gbl}^2 as:

$$q^*(\sigma_{\text{gbl}}^2) \sim \text{Inverse-Gamma} \left(\frac{1}{2}(L_{\text{gbl}} + 1), B_q(\sigma_{\text{gbl}}^2) \right),$$

where

$$B_q(\sigma_{\text{gbl}}^2) = \frac{1}{2} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u}^{\text{gbl}})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u}^{\text{gbl}})}) \right\} + \mu_{(1/a_{\text{gbl}})}$$

and using Result 1.4.3, we get

$$\mu_{q(1/\sigma_{\text{gbl}}^2)} = \frac{1}{2}(L_{\text{gbl}} + 1) / B_q(\sigma_{\text{gbl}}^2).$$

Expressions for $B_q(a_{\text{gbl}})$ and $\mu_{q(1/a_{\text{gbl}})}$

The derivations for $B_q(a_{\text{gbl}})$ and $\mu_{q(1/a_{\text{gbl}})}$ are similar to that for $B_q(\sigma_\varepsilon^2)$ and $\mu_{q(1/\sigma_\varepsilon^2)}$, respectively. Without going into further detail, the optimal q -density for a_{gbl} satisfies

$$q^*(a_{\text{gbl}}) \sim \text{Inverse-Gamma} \left(1, B_q(a_{\text{gbl}}) \right),$$

where

$$B_q(a_{\text{gbl}}) = \mu_{q(1/\sigma_{\text{gbl}}^2)} + A_{\text{gbl}}^{-2}$$

and

$$\mu_{q(1/a_{\text{gbl}})} = 1/B_{q(a_{\text{gbl}})}.$$

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

The full conditional distribution for Σ is given by:

$$\begin{aligned} p(\Sigma|\text{rest}) &\propto p(\mathbf{u}|\sigma_{\text{gbl}}^2, \Sigma, \sigma_{\text{grp}}^2) p(\Sigma|a_{\Sigma,1}, a_{\Sigma,2}) \\ &= |\mathbf{G}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{u}^\top \mathbf{G}^{-1}\mathbf{u}\right) |\Sigma|^{-\frac{(\nu+p-1)+p+1}{2}} \\ &\quad \times \exp\left[-\frac{1}{2}\text{tr}\{2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})\} \Sigma^{-1}\right] + \text{const} \end{aligned}$$

where $p = 2$ is the dimension of Σ . Taking the logarithm of both sides, we get:

$$\begin{aligned} \log p(\Sigma|\text{rest}) &= -\frac{m}{2} \log |\Sigma| - \frac{1}{2} \boldsymbol{\delta}^\top (\mathbf{I}_m \otimes \Sigma) \boldsymbol{\delta} - \frac{\nu+4}{2} \log |\Sigma| \\ &\quad - \text{tr}\{2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2}) \Sigma^{-1}\} + \text{const} \end{aligned}$$

where $\boldsymbol{\delta} = [\boldsymbol{\delta}_1^\top, \dots, \boldsymbol{\delta}_m^\top]^\top$. Further simplification results in

$$\begin{aligned} \log p(\Sigma|\text{rest}) &= -\frac{1}{2}(\nu + m + 4) \log |\Sigma| - \frac{1}{2} \text{tr}\left[\left\{\sum_{i=1}^m \boldsymbol{\delta}_i \boldsymbol{\delta}_i^\top + 2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})\right\} \Sigma^{-1}\right] \\ &\quad + \text{const}. \end{aligned}$$

Taking expectations:

$$\begin{aligned} \log q^*(\Sigma) &= E_q\{\log p(\Sigma|\text{rest})\} \\ &= -\frac{1}{2}(\nu + m + 4) \log |\Sigma| - \frac{1}{2} \left\{ \left[\sum_{i=1}^m \left\{ \boldsymbol{\mu}_{q(\boldsymbol{\delta}_i)} \boldsymbol{\mu}_{q(\boldsymbol{\delta}_i)}^\top + \Sigma_{q(\boldsymbol{\delta}_i)} \right\} \right. \right. \\ &\quad \left. \left. + 2\nu \text{diag}(\mu_{q(1/a_{\Sigma,1})}, \mu_{q(1/a_{\Sigma,2})}) \right] \Sigma^{-1} \right\} + \text{const} \end{aligned}$$

Therefore, we can represent the optimal q -density of Σ as:

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

where

$$B_{q(\Sigma)} = \sum_{i=1}^m \left\{ \boldsymbol{\mu}_{q(\boldsymbol{\delta}_i)} \boldsymbol{\mu}_{q(\boldsymbol{\delta}_i)}^\top + \Sigma_{q(\boldsymbol{\delta}_i)} \right\} + 2\nu \text{diag}(\mu_{q(1/a_{\Sigma,1})}, \mu_{q(1/a_{\Sigma,2})})$$

and using Results 1.4.7 and 1.4.8, we get

$$M_{q(\Sigma^{-1})} \equiv E_{q(\Sigma^{-1})} = (\nu + m + 1) B_{q(\Sigma)}^{-1}.$$

Expressions for the $B_{q(a_{\Sigma,j})}$ and $\mu_{q(1/a_{\Sigma,j})}$

We start with

$$\begin{aligned} p(a_{\Sigma,j} | \text{rest}) &\propto p(\boldsymbol{\Sigma} | a_{\Sigma,1}, a_{\Sigma,2}) p(a_{\Sigma,j}) \\ &= |2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})|^{\frac{\nu+1}{2}} \exp\left\{-\frac{1}{2} \text{tr}\left(2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2}) \boldsymbol{\Sigma}^{-1}\right)\right\} \\ &\quad \times (a_{\Sigma,j})^{-\frac{1}{2}-1} \exp\left\{-(1/A_{\Sigma,j})/a_{\Sigma,j}\right\} + \text{const} \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned} \log p(a_{\Sigma,j} | \text{rest}) &= -\frac{\nu+1}{2} \log(a_{\Sigma,j}) - \frac{3}{2} \log(a_{\Sigma,j}) - (1/A_{\Sigma,j})/a_{\Sigma,j} - \left\{\nu (\boldsymbol{\Sigma}^{-1})_{jj}\right\}/a_{\Sigma,j} \\ &\quad + \text{const} \\ &= -\frac{\nu+4}{2} \log(a_{\Sigma,j}) - \left\{\nu (\boldsymbol{\Sigma}^{-1})_{jj} + A_{\Sigma,j}^{-2}\right\}/a_{\Sigma,j} \end{aligned}$$

Taking expectations gives

$$\begin{aligned} \log q^*(a_{\Sigma,j}) &= E_q \{\log p(a_{\Sigma,j} | \text{rest})\} \\ &= -\frac{\nu+4}{2} \log(a_{\Sigma,j}) - \left\{\nu \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})}\right)_{jj} + A_{\Sigma,j}^{-2}\right\}/a_{\Sigma,j} + \text{const} \end{aligned}$$

Therefore,

$$q^*(a_{\Sigma,j}) \sim \text{Inverse-Gamma}\left(\frac{\nu}{2} + 1, B_{q(a_{\Sigma,j})}\right), \quad j = 1, 2,$$

where

$$B_{q(a_{\Sigma,j})} = \nu \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})}\right)_{jj} + A_{\Sigma,j}^{-2}$$

and through the use of Result 1.4.3

$$\mu_{q(1/a_{\Sigma,j})} = \left(\frac{\nu}{2} + 1\right) / B_{q(a_{\Sigma,j})}.$$

Expressions for $B_{q(\sigma_{\text{grp}}^2)}$ and $\mu_{q(1/\sigma_{\text{grp}}^2)}$

The derivations for $B_{q(\sigma_{\text{grp}}^2)}$ and $\mu_{q(1/\sigma_{\text{grp}}^2)}$ are similar to that for $B_{q(\sigma_{\text{gbl}}^2)}$ and $\mu_{q(1/\sigma_{\text{gbl}}^2)}$, respectively. Without going through further detail, the optimal q -density for σ_{grp}^2 is shown to be of the form:

$$q^*(\sigma_{\text{grp}}^2) \sim \text{Inverse-Gamma}\left(\frac{1}{2}(m L_{\text{grp}} + 1), B_{q(\sigma_{\text{grp}}^2)}\right),$$

where

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

and using Result 1.4.3, we get

$$\mu_{q(1/\sigma_{\text{grp}}^2)} = \frac{1}{2}(m L_{\text{grp}} + 1)/B_{q(\sigma_{\text{grp}}^2)}.$$

5.B Derivation of the marginal log-likelihood lower bound

The marginal log-likelihood lower bound expression given in (5.13) is:

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q) &= \frac{1}{2}(\nu + p - 1) \log(2\nu) - \frac{1}{2} \sum_{i=1}^m n_i \log(2\pi) - 4 \log(\pi) - \frac{p}{2} \log(\sigma_\beta^2) \\ &\quad - \frac{1}{\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| \\ &\quad + \frac{1}{2} \{p + L_{\text{gbl}} + m(p + L_{\text{grp}})\} - \log(C_{p, \nu+p-1}) + \log(C_{p, \nu+m+p-1}) \\ &\quad - \frac{1}{2}(\nu + p + m - 1) \log |\mathbf{B}_{q(\boldsymbol{\Sigma})}| + \log \Gamma \left\{ \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \right\} \\ &\quad - \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \log(B_{q(\sigma_\varepsilon^2)}) + \log \Gamma \left\{ \frac{1}{2} (L_{\text{gbl}} + 1) \right\} \\ &\quad - \frac{1}{2} (L_{\text{gbl}} + 1) \log(B_{q(\sigma_{\text{gbl}}^2)}) + \log \Gamma \left\{ \frac{1}{2} (m \times L_{\text{grp}} + 1) \right\} \\ &\quad - \frac{1}{2} (m \times L_{\text{grp}} + 1) \log(B_{q(\sigma_{\text{grp}}^2)}) - \log(A_\varepsilon) - \log(B_{q(a_\varepsilon)}) \\ &\quad + \mu_{q(1/\sigma_\varepsilon^2)} \mu_{q(1/a_\varepsilon)} - \log(B_{q(a_{\text{grp}})}) + \sum_{j=1}^p \nu \left(\mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma, j})} \\ &\quad - \log(A_{\text{gbl}}) - \log(B_{q(a_{\text{gbl}})}) + \mu_{q(1/\sigma_{\text{gbl}}^2)} \mu_{q(1/a_{\text{gbl}})} - \log(A_{\text{grp}}) \\ &\quad + \mu_{q(1/\sigma_{\text{grp}}^2)} \mu_{q(1/\sigma_{\text{grp}}^2)} - \sum_{j=1}^p \log(A_{\Sigma, j}) + p \log \Gamma \left\{ \frac{1}{2} (\nu + p) \right\} \\ &\quad - \frac{1}{2} (\nu + p - 1) \sum_{j=1}^p \log(B_{q(a_{\Sigma, j})}) \end{aligned}$$

5.B. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

The details of the derivation are as follows:

$$\begin{aligned}
\log \underline{p}(\mathbf{y}; q) &= E_q \left\{ \log p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2, a_\varepsilon, a_{\text{gbl}}, \mathbf{a}_\Sigma, a_{\text{grp}}) \right. \\
&\quad \left. - \log q^*(\boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2, \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2, a_\varepsilon, a_{\text{gbl}}, \mathbf{a}_\Sigma, a_{\text{grp}}) \right\} \\
&= E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \right. \\
&\quad + E_q \left\{ \log p(\boldsymbol{\beta}, \mathbf{u} | \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}) \right\} \\
&\quad + E_q \left\{ \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \right\} \\
&\quad + E_q \left\{ \log p(\sigma_{\text{gbl}}^2 | a_{\text{gbl}}) - \log q^*(\sigma_{\text{gbl}}^2) \right\} \\
&\quad + E_q \left\{ \log p(\boldsymbol{\Sigma} | a_{\Sigma,1}, a_{\Sigma,2}) - \log q^*(\boldsymbol{\Sigma}) \right\} \\
&\quad + E_q \left\{ \log p(\sigma_{\text{grp}}^2 | a_{\text{grp}}) - \log q^*(\sigma_{\text{grp}}^2) \right\} \\
&\quad + E_q \left\{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \right\} + E_q \left\{ \log p(a_{\text{gbl}}) - \log q^*(a_{\text{gbl}}) \right\} \\
&\quad \left. + E_q \left\{ \log p(\mathbf{a}_\Sigma) - \log q^*(\mathbf{a}_\Sigma) \right\} + E_q \left\{ \log p(a_{\text{grp}}) - \log q^*(a_{\text{grp}}) \right\} \right\}.
\end{aligned} \tag{5.23}$$

First note that

$$\log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) = \frac{1}{2} \sum_{i=1}^m n_i \log(2\pi) - \frac{1}{2} \sum_{i=1}^m n_i \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2$$

Taking expectations we get

$$\begin{aligned}
E_q \left\{ \log p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2) \right\} &= -\frac{1}{2} \sum_{i=1}^m n_i \log(2\pi) - \frac{1}{2} \sum_{i=1}^m n_i E_q \left\{ \log(\sigma_\varepsilon^2) \right\} \\
&\quad - \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} E_q \left(\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2 \right).
\end{aligned} \tag{5.24}$$

Next, using the concatenated form of the design matrices, $\mathbf{C} = [\mathbf{X} | \mathbf{Z}]$ and Results 1.4.19 and 1.4.12, we see that

$$\begin{aligned}
E_q \left(\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{u}\|^2 \right) &= E_q \left(\left\| \mathbf{y} - \mathbf{C} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right\|^2 \right) \\
&= \left\| \mathbf{y} - \mathbf{C} E_q \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right) \right\|^2 + \text{tr} \left\{ \mathbf{C} \text{Cov} \left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \right) \mathbf{C}^\top \right\} \\
&= \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^\top) \\
&= \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}).
\end{aligned}$$

Substituting into (5.24), we have

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

5.B. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

Moving on to the next term in (5.23), we see that

$$\begin{aligned}
\log p(\boldsymbol{\beta}, \mathbf{u} | \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2) &= -\log(2\pi) - \log(\sigma_{\beta}^2) - \frac{1}{2\sigma_{\beta}^2} \|\boldsymbol{\beta}\|^2 \\
&\quad - \frac{1}{2} \{L_{\text{gbl}} + m(2 + L_{\text{grp}})\} \log(2\pi) - \frac{L_{\text{gbl}}}{2} \log(\sigma_{\text{gbl}}^2) \\
&\quad - \frac{m}{2} \log|\boldsymbol{\Sigma}| - \frac{1}{2} m \times L_{\text{grp}} \log(\sigma_{\text{grp}}^2) - \frac{1}{2\sigma_{\text{gbl}}^2} \|\mathbf{u}^{\text{gbl}}\|^2 \\
&\quad - \frac{1}{2} \text{tr} \left(\sum_{i=1}^m \boldsymbol{\delta}_i \boldsymbol{\delta}_i^{\top} \boldsymbol{\Sigma}^{-1} \right) - \frac{1}{2\sigma_{\text{grp}}^2} \|\mathbf{u}^{\text{grp}}\|^2.
\end{aligned}$$

Next, note that the entropy of the multivariate normal random vector $[\boldsymbol{\beta}^{\top}, \mathbf{u}^{\top}]^{\top}$ is

$$\begin{aligned}
-E_q \{ \log q^*(\boldsymbol{\beta}, \mathbf{u}) \} &= \frac{1}{2} \log \left\{ (2\pi e)^{2+L_{\text{gbl}}+m(2+L_{\text{grp}})} |\boldsymbol{\Sigma}_q(\boldsymbol{\beta}, \mathbf{u})| \right\} \\
&= \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\beta}, \mathbf{u})| + \frac{2+L_{\text{gbl}}+m(2+L_{\text{grp}})}{2} \log(2\pi) + \frac{2+L_{\text{gbl}}+m(2+L_{\text{grp}})}{2}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
E_q \{ \log p(\boldsymbol{\beta}, \mathbf{u} | \sigma_{\text{gbl}}^2, \boldsymbol{\Sigma}, \sigma_{\text{grp}}^2) - \log q^*(\boldsymbol{\beta}, \mathbf{u}) \} \\
&= -\log(2\pi) - \log(\sigma_{\beta}^2) - \frac{1}{2\sigma_{\beta}^2} \{ \|\boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\beta})) \} \\
&\quad - \frac{1}{2} \{L_{\text{gbl}} + m(2 + L_{\text{grp}})\} \log(2\pi) - \frac{L_{\text{gbl}}}{2} E_q \{ \log(\sigma_{\text{gbl}}^2) \} \\
&\quad - \frac{m}{2} E_q \{ \log|\boldsymbol{\Sigma}| \} - \frac{1}{2} m L_{\text{grp}} E_q \{ \log(\sigma_{\text{grp}}^2) \} \\
&\quad - \frac{1}{2} \mu_q(1/\sigma_{\text{gbl}}^2) \{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{gbl}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{gbl}})) \} \\
&\quad - \frac{1}{2} \text{tr} \left\{ \sum_{i=1}^m \left(\boldsymbol{\mu}_q(\boldsymbol{\delta}_i) \boldsymbol{\mu}_q^{\top}(\boldsymbol{\delta}_i) + \boldsymbol{\Sigma}_q(\boldsymbol{\delta}_i) \right) \mathbf{M}_q(\boldsymbol{\Sigma}^{-1}) \right\} \\
&\quad - \frac{1}{2} \mu_q(1/\sigma_{\text{grp}}^2) \{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{grp}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{grp}})) \} \\
&\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\beta}, \mathbf{u})| + \frac{1}{2} (2 + L_{\text{gbl}} + m(2 + L_{\text{grp}})) \log(2\pi) \\
&\quad + \frac{1}{2} (2 + L_{\text{gbl}} + m(2 + L_{\text{grp}})).
\end{aligned}$$

Collecting like terms and removing cancellations, this equates to

$$\begin{aligned}
&= -\log(\sigma_{\beta}^2) - \frac{1}{2\sigma_{\beta}^2} \{ \|\boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\boldsymbol{\beta})) \} - \frac{L_{\text{gbl}}}{2} E_q \{ \log(\sigma_{\text{gbl}}^2) \} \\
&\quad - \frac{m}{2} E_q \{ \log|\boldsymbol{\Sigma}| \} - \frac{1}{2} m L_{\text{grp}} E_q \{ \log(\sigma_{\text{grp}}^2) \} \\
&\quad - \frac{1}{2} \mu_q(1/\sigma_{\text{gbl}}^2) \{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{gbl}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{gbl}})) \} \\
&\quad - \frac{1}{2} \text{tr} \left\{ \sum_{i=1}^m \left(\boldsymbol{\mu}_q(\boldsymbol{\delta}_i) \boldsymbol{\mu}_q^{\top}(\boldsymbol{\delta}_i) + \boldsymbol{\Sigma}_q(\boldsymbol{\delta}_i) \right) \mathbf{M}_q(\boldsymbol{\Sigma}^{-1}) \right\} \\
&\quad - \frac{1}{2} \mu_q(1/\sigma_{\text{grp}}^2) \{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{grp}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{grp}})) \} \\
&\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_q(\boldsymbol{\beta}, \mathbf{u})| + \frac{1}{2} (2 + L_{\text{gbl}} + m(2 + L_{\text{grp}})).
\end{aligned}$$

Next,

$$\begin{aligned}
 \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) &= \\
 & \log \left\{ \frac{(1/a_\varepsilon)^{\frac{1}{2}}}{\Gamma(\frac{1}{2})} (\sigma_\varepsilon^2)^{-\frac{1}{2}-1} \exp\left(-\frac{1/a_\varepsilon}{\sigma_\varepsilon^2}\right) \right\} \\
 & - \log \left\{ \frac{B_q(\sigma_\varepsilon^2)^{\frac{1}{2}(\sum_{i=1}^m n_i + 1)}}{\Gamma(\frac{1}{2}(\sum_{i=1}^m n_i + 1))} (\sigma_\varepsilon^2)^{-\frac{1}{2}(\sum_{i=1}^m n_i + 1)-1} \exp\left(-\frac{B_q(\sigma_\varepsilon^2)}{\sigma_\varepsilon^2}\right) \right\} \\
 &= -\frac{1}{2} \log(a_\varepsilon) - \frac{3}{2} \log(\sigma_\varepsilon^2) - (1/a_\varepsilon)/\sigma_\varepsilon^2 - \frac{1}{2} \log(\pi) \\
 & - \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) \log(B_q(\sigma_\varepsilon^2)) + \left\{ \frac{1}{2} \left(\sum_{i=1}^m n_i + 1 \right) + 1 \right\} \log(\sigma_\varepsilon^2) \\
 & + B_q(\sigma_\varepsilon^2)/\sigma_\varepsilon^2 + \log \left\{ \Gamma\left(\frac{1}{2} \sum_{i=1}^m n_i + 1\right) \right\}
 \end{aligned}$$

Thus,

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

The derivation for the term

$$E_q \left\{ \log p(\sigma_{\text{gbl}}^2 | a_{\text{gbl}}) - \log q^*(\sigma_{\text{gbl}}^2) \right\}$$

is similar to that for

$$E_q \left\{ \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \right\}.$$

So, without going into further detail,

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

We now move on to expanding the fifth term in (5.23):

$$\begin{aligned}
 \log p(\boldsymbol{\Sigma} | a_{\Sigma,1}, a_{\Sigma,2}) - \log q^*(\boldsymbol{\Sigma}) &= -\log(C_{2,\nu+2-1}) + \frac{\nu+2-1}{2} \log |2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})| \\
 & - \frac{\nu+4}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr} \left\{ 2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2}) \boldsymbol{\Sigma}^{-1} \right\} \\
 & + \log(C_{2,\nu+m+1}) - \frac{\nu+m+1}{2} \log |\mathbf{B}_q(\boldsymbol{\Sigma})| + \frac{\nu+4+m}{2} \log |\boldsymbol{\Sigma}| \\
 & + \frac{1}{2} \text{tr} \left\{ \mathbf{B}_q(\boldsymbol{\Sigma}) \boldsymbol{\Sigma}^{-1} \right\}
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 E_q \{ \log p(\boldsymbol{\Sigma} | a_{\Sigma,1}, a_{\Sigma,2}) - \log q^*(\boldsymbol{\Sigma}) \} = & \\
 & - \log(C_{2,\nu+1}) + \log(C_{2,\nu+m+1}) \\
 & + \frac{\nu+1}{2} E_q \{ \log |2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})| \} \\
 & - \frac{\nu+m+1}{2} \log |\mathbf{B}_{q(\boldsymbol{\Sigma})}| + \frac{m}{2} E_q \{ \log |\boldsymbol{\Sigma}| \} \\
 & + \frac{1}{2} \text{tr} \left\{ \left(\mathbf{B}_{q(\boldsymbol{\Sigma})} - 2\nu \text{diag} \left(\mu_{q(1/a_{\Sigma,1})}, \mu_{q(1/a_{\Sigma,2})} \right) \right) \mathbf{M}_{q(\boldsymbol{\Sigma}^{-1})} \right\}.
 \end{aligned}$$

The derivation for the term

$$E_q \{ \log p(\sigma_{\text{grp}}^2 | a_{\text{grp}}) - \log q^*(\sigma_{\text{grp}}^2) \}$$

is similar to that for

$$E_q \{ \log p(\sigma_\varepsilon^2 | a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \}.$$

Without going into further detail, it is shown that

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

Next, we have

$$\begin{aligned}
 \log p(a_\varepsilon) - \log q^*(a_\varepsilon) = & \log \left\{ \frac{(1/A_\varepsilon^2)^{\frac{1}{2}}}{\Gamma(\frac{1}{2})} (a_\varepsilon)^{-\frac{1}{2}-1} \exp \left(-\frac{1/A_\varepsilon^2}{a_\varepsilon} \right) \right\} \\
 & - \log \left\{ \frac{B_{q(a_\varepsilon)}}{\Gamma(1)} (a_\varepsilon)^{-1-1} \exp \left(-\frac{B_{q(a_\varepsilon)}}{a_\varepsilon} \right) \right\} \\
 = & -\log(A_\varepsilon) - \frac{1}{2} \log(\pi) - \frac{3}{2} \log(a_\varepsilon) - (1/A_\varepsilon^2)/a_\varepsilon - \log(B_{q(a_\varepsilon)}) \\
 & + 2 \log(a_\varepsilon) + B_{q(a_\varepsilon)}/a_\varepsilon.
 \end{aligned}$$

From algorithm 5, we see that

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

so making use of this and taking expectations we get:

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

5.B. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

The derivation of the following term in (5.23) is similar to that just shown above, thus we have

$$E_q \{ \log p(a_{\text{gbl}}) - \log q^*(a_{\text{gbl}}) \} = -\log(A_{\text{gbl}}) - \frac{1}{2} \log(\pi) + \frac{1}{2} E_q \{ \log(a_{\text{gbl}}) \} - \log(B_{q(a_{\text{gbl}})}) \\ + \mu_{q(1/\sigma_{\text{gbl}}^2)} \mu_{q(1/a_{\text{gbl}})}.$$

Next,

$$\log p(\mathbf{a}_\Sigma) - \log q^*(\mathbf{a}_\Sigma) = \log \left\{ \prod_{j=1}^2 \frac{(1/A_{\Sigma,j}^2)^{\frac{1}{2}}}{\Gamma(\frac{1}{2})} (a_{\Sigma,j})^{-\frac{1}{2}-1} \exp\left(-\frac{1/A_{\Sigma,j}^2}{a_{\Sigma,j}}\right) \right\} \\ - \log \left\{ \prod_{j=1}^2 \frac{B_{q(a_{\Sigma,j})}^{\frac{\nu}{2}+1}}{\Gamma(\frac{\nu}{2}+1)} (a_{\Sigma,j})^{-(\frac{\nu}{2}+1)-1} \exp\left(-\frac{B_{q(a_{\Sigma,j})}}{a_{\Sigma,j}}\right) \right\} \\ = -\sum_{j=1}^2 \log(A_{\Sigma,j}) - \log(\pi) + 2 \log\{\Gamma(\frac{\nu}{2}+1)\} + \frac{1}{2}(\nu+1) \sum_{j=1}^2 \log(a_{\Sigma,j}) \\ - \sum_{j=1}^2 (1/A_{\Sigma,j}^2)/a_{\Sigma,j} + \sum_{j=1}^2 B_{q(a_{\Sigma,j})}/a_{\Sigma,j} - (\frac{\nu}{2}+1) \sum_{j=1}^2 \log(B_{q(a_{\Sigma,j})}).$$

Using the fact that

$$B_{q(a_{\Sigma,j})} = \nu \left(\mathbf{M}_{q(\Sigma^{-1})} \right)_{jj} + 1/A_{\Sigma,j}^2,$$

and taking expectations, we get

$$E_q \{ \log p(\mathbf{a}_\Sigma) - \log q^*(\mathbf{a}_\Sigma) \} = -\sum_{j=1}^2 \log(A_{\Sigma,j}) - \log(\pi) + 2 \log\{\Gamma(\frac{\nu}{2}+1)\} \\ + \frac{1}{2}(\nu+1) \sum_{j=1}^2 E_q \{ \log(a_{\Sigma,j}) \} \\ - (\frac{\nu}{2}+1) \sum_{j=1}^2 \log(B_{q(a_{\Sigma,j})}) \\ + \sum_{j=1}^2 \nu \left(\mathbf{M}_{q(\Sigma^{-1})} \right)_{jj} \mu_{q(1/a_{\Sigma,j})}.$$

The derivation for the last term in (5.23) is similar to the derivation for the a_ε parameter.

Thus, the expression is:

$$E_q \{ \log p(a_{\text{grp}}) - \log q^*(a_{\text{grp}}) \} = -\log(A_{\text{grp}}) - \frac{1}{2} \log(\pi) + \frac{1}{2} E_q \{ \log(a_{\text{grp}}) \} \\ - \log(B_{q(a_{\text{grp}})}) + \mu_{q(1/\sigma_{\text{grp}}^2)} \mu_{q(1/a_{\text{grp}})}$$

Bringing all these terms together, and noting the cancellations:

$$-\frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} \left\{ \|\mathbf{y} - \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})}\|^2 + \text{tr}(\mathbf{C}^\top \mathbf{C} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) \right\} = (B_{q(\sigma_\varepsilon^2)} - \mu_{q(1/a_\varepsilon)}) \mu_{q(1/\sigma_\varepsilon^2)},$$

5.B. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

$$\begin{aligned}
& -\frac{1}{2}\mu_q(1/\sigma_{\text{gbl}}^2) \left\{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{gbl}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{gbl}})) \right\} = \left(B_q(\sigma_{\text{gbl}}^2) - \mu_q(1/a_{\text{gbl}}) \right) \mu_q(1/\sigma_{\text{gbl}}^2), \\
& -\frac{1}{2}\text{tr} \left\{ \sum_{i=1}^m \left(\boldsymbol{\mu}_q(\boldsymbol{\delta}_i) \boldsymbol{\mu}_q^\top(\boldsymbol{\delta}_i) + \boldsymbol{\Sigma}_q(\boldsymbol{\delta}_i) \right) \mathbf{M}_q(\boldsymbol{\Sigma}^{-1}) \right\} \\
& \quad = \frac{1}{2}\text{tr} \left\{ \left(\mathbf{B}_q(\boldsymbol{\Sigma}) - 2\nu \text{diag}(\mu_q(1/a_{\Sigma,1}), \mu_q(1/a_{\Sigma,2})) \right) \mathbf{M}_q(\boldsymbol{\Sigma}^{-1}) \right\}, \\
& -\frac{1}{2}\mu_q(1/\sigma_{\text{grp}}^2) \left\{ \|\boldsymbol{\mu}_q(\mathbf{u}^{\text{grp}})\|^2 + \text{tr}(\boldsymbol{\Sigma}_q(\mathbf{u}^{\text{grp}})) \right\} = \left(B_q(\sigma_{\text{grp}}^2) - \mu_q(1/a_{\text{grp}}) \right) \mu_q(1/\sigma_{\text{grp}}^2)
\end{aligned}$$

and the fact that

$$\frac{\nu+1}{2} E_q \{ \log |2\nu \text{diag}(1/a_{\Sigma,1}, 1/a_{\Sigma,2})| \} + \frac{\nu+1}{2} \sum_{j=1}^2 E_q \{ \log(a_{\Sigma,j}) \} = \frac{\nu+1}{2} \log(2\nu),$$

gives the final lower bound expression in (5.13) where $p = 2$ is the column length in \mathbf{X} .

Chapter 6

Mean field variational Bayes for mixture models in measurement error problems

6.1 Introduction

In the area of epidemiology it is generally recognised that covariates, such as risk factors, are not being measured accurately on all subjects. This can heavily affect the estimate of the relationship between these covariates and the outcome. These are known as *measurement error* or *errors-in-variables* problems. In this chapter we examine the Bayesian formulation of such a problem, in particular fast mean field variational Bayes approximation and compare this to the MCMC benchmark. The field of epidemiology has lent great motivation to this chapter, as measurement error has long been a problem there. The model we present, inspired by the model used in Richardson *et al.* (2002), can be applied in a variety of settings involving measurement error.

The main concern for problems with measurement error is to make inference on the relationship between a response \mathbf{y} and covariates \mathbf{x} , in situations where the only other information we have on \mathbf{x} is through the recording of an error contaminated surrogate \mathbf{o} , as \mathbf{x} may not have been measured accurately on all subjects. An attempt to regress \mathbf{y} on \mathbf{o} whilst ignoring the problem of measurement error can in most cases lead to extremely inaccurate inferential statistics. An extensive literature has formed for flexible estimation of regression models where precise measurements on the predictor/s are not available. An

extensive reference to the area is Carroll *et al.* (2006). The 2000s in particular saw an influx of research into the area, e.g. Mallick *et al.* (2002), Liang *et al.* (2003), Carroll *et al.* (2004), Ganguli *et al.* (2005) and Carroll *et al.* (2007). Many contributions have discussed the use of MCMC based inference to fitting models impaired by measurement error, however inference based on MCMC can at times be extremely time intensive. As an alternative to MCMC, we focus on MFVB fitting to the model structure given in Richardson *et al.* (2002) where the response follows a normal distribution.

Section 6.2 discusses the different components necessary for the measurement error model with normal mixture components. In Section 6.3 we present the full measurement error model that we deal with. We present the MFVB algorithm in Section 6.4 and Section 6.5 provides a simulation study based on the new MFVB algorithm. The appendices provide details on the calculations required for the MFVB algorithm and corresponding lower bound expression.

6.2 Model structure

We set up the covariate of interest \mathbf{x} by splitting it up into two subgroups: \mathbf{x}_{obs} which involves the accurate measurements of the covariate; and $\mathbf{x}_{\text{unobs}}$ which involves the part of the covariate that has not been measured accurately. The observations in \mathbf{x}_{obs} are generally recorded using a *gold standard method*, which is an error-free method for measuring the covariate, however this is costly to use on a large scale, hence $\mathbf{x}_{\text{unobs}}$ is usually much larger than \mathbf{x}_{obs} . The subjects belonging to \mathbf{x}_{obs} are usually referred to as the *validation group*, and the subjects belonging to $\mathbf{x}_{\text{unobs}}$ as the *main study*. This is in keeping with the terminology used in Richardson *et al.* (2002).

Throughout this chapter we let i refer to a particular individual, where $1 \leq i \leq n$. Also \mathbf{y} denotes the known outcome (e.g. disease status), \mathbf{x} the true covariate comprising observed quantities \mathbf{x}_{obs} and unobserved quantities $\mathbf{x}_{\text{unobs}}$, \mathbf{o} the observed surrogate for \mathbf{x} and \mathbf{c} denotes a known covariate. We focus on the case where \mathbf{x} and \mathbf{o} are univariate. The general model is divided into three submodels:

1. a *regression model* expressing the relationship between the covariates \mathbf{c} and \mathbf{x} and the outcome \mathbf{y} , denoted by $p(\mathbf{y}|\mathbf{x},\boldsymbol{\beta})$;
2. a *measurement model*, which expresses the relationship between the surrogate mea-

sures \mathbf{o} and the true unknown covariate \mathbf{x} , denoted by $p(\mathbf{o}|\mathbf{x}, \boldsymbol{\alpha})$; and

3. the *prior model* for \mathbf{x} , which specifies the distribution of the unknown covariate \mathbf{x} in the general population and is denoted by $p(\mathbf{x}|\boldsymbol{\theta})$,

where $\boldsymbol{\beta}$ is the regression coefficient vector for \mathbf{y} , $\boldsymbol{\alpha}$ is the regression coefficient vector for \mathbf{o} and $\boldsymbol{\theta}$ are the parameters in the prior model for \mathbf{x} .

6.2.1 Nondifferential measurement error

An important assumption for our model is one of nondifferential measurement error which holds when \mathbf{o} has no information about \mathbf{y} other than what is given in \mathbf{x} and \mathbf{c} . Thus, if \mathbf{o} is conditionally independent of \mathbf{y} given \mathbf{x} and \mathbf{c} , then it is a surrogate. This can be expressed as

$$p(\mathbf{y}|\mathbf{o}, \mathbf{x}, \mathbf{c}, \boldsymbol{\beta}) = p(\mathbf{y}|\mathbf{x}, \mathbf{c}, \boldsymbol{\beta}) \quad \text{or} \quad \mathbf{y} \perp \mathbf{o}|\mathbf{x}.$$

This leads to the joint distribution

$$p(\boldsymbol{\theta}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{x}, \mathbf{o}, \mathbf{y}) = p(\boldsymbol{\theta})p(\boldsymbol{\alpha})p(\boldsymbol{\beta}) \prod_{i=1}^n p(x_i|\boldsymbol{\theta}) \prod_{i=1}^n p(o_i|x_i, \boldsymbol{\alpha}) \prod_{i=1}^n p(y_i|x_i, c_i, \boldsymbol{\beta}), \quad (6.1)$$

where $\boldsymbol{\theta}$ denotes all the parameters of the normal mixture model and $p(\boldsymbol{\beta})$, $p(\boldsymbol{\alpha})$ and $p(\boldsymbol{\theta})$ are the prior distributions for the parameters in the three submodels.

6.2.2 Finite normal mixture component

A question that has been heavily researched in the last decade is that of how to model $p(\mathbf{x}|\boldsymbol{\theta})$ in such that we allow the data to influence its shape, since regression parameters are known to be susceptible to the particular shape of $p(\mathbf{x}|\boldsymbol{\theta})$. Many contributions have been made to this particular question in the measurement error context, some of which are described in Carroll *et al.* (1993), Roeder *et al.* (1996), Schafer (2001), Aitkin & Rocci (2002) and Müller & Roeder (1997). Richardson *et al.* (2002) however, suggest mixture models with alterable number of components as a natural alternative. In particular, mixtures of normal distributions, as they make model misspecification for $p(\mathbf{x}|\boldsymbol{\theta})$ increasingly more robust. Thus, we also impose a mixture of normal distributions on $p(\mathbf{x}|\boldsymbol{\theta})$ of the form:

$$p(x_i|a_{i1}^x, \dots, a_{iK}^x) = \prod_{k=1}^K \left[(2\pi\sigma_k^2)^{-1/2} \exp \left\{ -\frac{1}{2}(x_i - \mu_k)^2 / \sigma_k^2 \right\} \right]^{a_{ik}^x}, \quad (6.2)$$

$$(a_{i1}^x, \dots, a_{iK}^x) | (\omega_1, \dots, \omega_K) \stackrel{\text{ind.}}{\sim} \text{Multinomial}(1; \omega_1, \dots, \omega_K), \quad 1 \leq i \leq n,$$

where the a_{ik} are auxiliary variables introduced to make the mixture model tractable and K refers to the number of mixture components used. An example of this use of auxiliary variables can be found in Section 2.2.4 of Ormerod & Wand (2010). The prior distributions have the following forms:

$$\begin{aligned} \mu_k &\stackrel{\text{ind.}}{\sim} \text{N}(\mu_{\mu_k}, \sigma_{\mu_k}^2), \quad \sigma_k \stackrel{\text{ind.}}{\sim} \text{Half-Cauchy}(A_k), \quad 1 \leq k \leq K, \\ (\omega_1, \dots, \omega_K) &\sim \text{Dirichlet}(\alpha, \dots, \alpha), \quad \alpha > 0, \end{aligned} \tag{6.3}$$

where according to this notation, $\sum_{k=1}^K a_{ik}^x = 1$, indicating which mixture component the current observation belongs to, and $\omega_k = P(a_{ik}^x = 1)$.

6.2.3 Joint model

Bringing the distributions in (6.1), (6.2) and (6.3) together and assuming natural conditional independencies, our joint model for the measurement error problem is expressed as

$$\begin{aligned} &\prod_{i=1}^n \{p(y_i | x_i, \mathbf{c}_i, \boldsymbol{\beta})\} \prod_{i=1}^n \{p(o_i | x_i, \boldsymbol{\alpha})\} \prod_{i=1}^n \prod_{k=1}^K \{p(x_i | a_{ik}, \mu_k, \sigma_k^2)\} \\ &\times \prod_{i=1}^n \{p(a_{i1}^x, \dots, a_{iK}^x)\} p(\omega_1, \dots, \omega_K) \prod_{k=1}^K \{p(\mu_k) p(\sigma_k)\} p(\boldsymbol{\beta}) p(\boldsymbol{\alpha}). \end{aligned}$$

The model components discussed in the previous sections are brought together next to present the full model, considered in a measurement error setting.

6.3 The full model

We consider the model given in (6.4) which exploits the structure described in Section 6.2.

$$\begin{aligned}
 & \left. \begin{aligned}
 y_i | x_i, \boldsymbol{\beta}, \sigma_\varepsilon^2 &\stackrel{\text{ind.}}{\sim} \text{N}\{(\mathbf{X}\boldsymbol{\beta})_i, \sigma_\varepsilon^2\}, \quad \sigma_\varepsilon^2 | a_\varepsilon \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_\varepsilon\right), \\
 a_\varepsilon &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_\varepsilon^2\right), \\
 \boldsymbol{\beta} &\sim \text{N}(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_p),
 \end{aligned} \right\} \begin{array}{l} \text{regression} \\ \text{model} \end{array} \\
 & \left. \begin{aligned}
 o_i | x_i, \boldsymbol{\alpha}, \sigma_o^2 &\stackrel{\text{ind.}}{\sim} \text{N}\{(\tilde{\mathbf{X}}\boldsymbol{\alpha})_i, \sigma_o^2\}, \quad \boldsymbol{\alpha} \sim \text{N}(\mathbf{0}, \sigma_\alpha^2 \mathbf{I}_q), \\
 \sigma_o^2 | a_o &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_o\right), \quad a_o \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_o^2\right),
 \end{aligned} \right\} \begin{array}{l} \text{measurement} \\ \text{error} \\ \text{model} \end{array} \quad (6.4) \\
 & \left. \begin{aligned}
 x_i | \mathbf{a}_i, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 &= \prod_{k=1}^K \left[\{2\pi (\sigma_k^x)^2\}^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \mu_k^x)^2 / (\sigma_k^x)^2\right\} \right]^{a_{ik}}, \\
 a_{i1}, \dots, a_{iK} &\sim \text{Multinomial}(1; \omega_1, \dots, \omega_K), \quad \omega_1, \dots, \omega_K \sim \text{Dirichlet}(\alpha, \dots, \alpha), \\
 \mu_k^x &\sim \text{N}(\mu_\mu, \sigma_\mu^2), \quad (\sigma_k^x)^2 | a_k^x \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_k^x\right), \\
 a_k^x &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/(A_k^x)^2\right), \quad 1 \leq k \leq K,
 \end{aligned} \right\} \begin{array}{l} \text{prior} \\ \text{specification} \end{array}
 \end{aligned}$$

where \mathbf{a} is an $n \times K$ matrix of auxiliary variables and \mathbf{a}_i is the i th column of \mathbf{a} with entries a_{i1}, \dots, a_{iK} . In addition, $\boldsymbol{\mu}^x \equiv (\mu_1^x, \dots, \mu_K^x)$ and $(\boldsymbol{\sigma}^x)^2 \equiv \{(\sigma_1^x)^2, \dots, (\sigma_K^x)^2\}$. We also define

$$\boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0 \\ \beta_{\mathbf{c}} \\ \beta_{\mathbf{x}} \end{bmatrix}, \quad \boldsymbol{\alpha} \equiv \begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} 1 & c_1 & x_1 \\ \vdots & \vdots & \vdots \\ 1 & c_n & x_n \end{bmatrix}, \quad \tilde{\mathbf{X}} \equiv \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}.$$

Here $\mathbf{y} = (y_1, \dots, y_n)$ is an $n \times 1$ vector of response variables, $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$ are the $p \times 1$ and $q \times 1$ vectors of fixed effects, where $p = 3$ and $q = 2$, \mathbf{X} and $\tilde{\mathbf{X}}$ are corresponding design matrices. Also, σ_ε^2 and σ_o^2 are the variance parameters corresponding to the fixed effects vector and a_ε and a_o are the corresponding auxiliary variables as represented in Result 1.4.6. The variables in the *prior specification* of (6.4) are as described in Section 6.2.2.

We define the following additional notation which is useful for the remaining sections of this chapter. Let n_{obs} denote the number of observed x_i measurements and n_{unobs} denote the number of unobserved x_i measurements. We define \mathbf{x}_{obs} to be the $n_{\text{obs}} \times 1$ vector containing the observed x_i measurements and $\mathbf{x}_{\text{unobs}}$ to be the $n_{\text{unobs}} \times 1$ vector containing the unobserved x_i measurements. In addition, let $\mathbf{y}_{x_{\text{obs}}}$ and $\mathbf{y}_{x_{\text{unobs}}}$ be the $n_{\text{obs}} \times 1$ and $n_{\text{unobs}} \times 1$ vectors corresponding to the measurements in \mathbf{x}_{obs} and $\mathbf{x}_{\text{unobs}}$. $\mathbf{c}_{x_{\text{obs}}}$, and $\mathbf{c}_{x_{\text{unobs}}}$

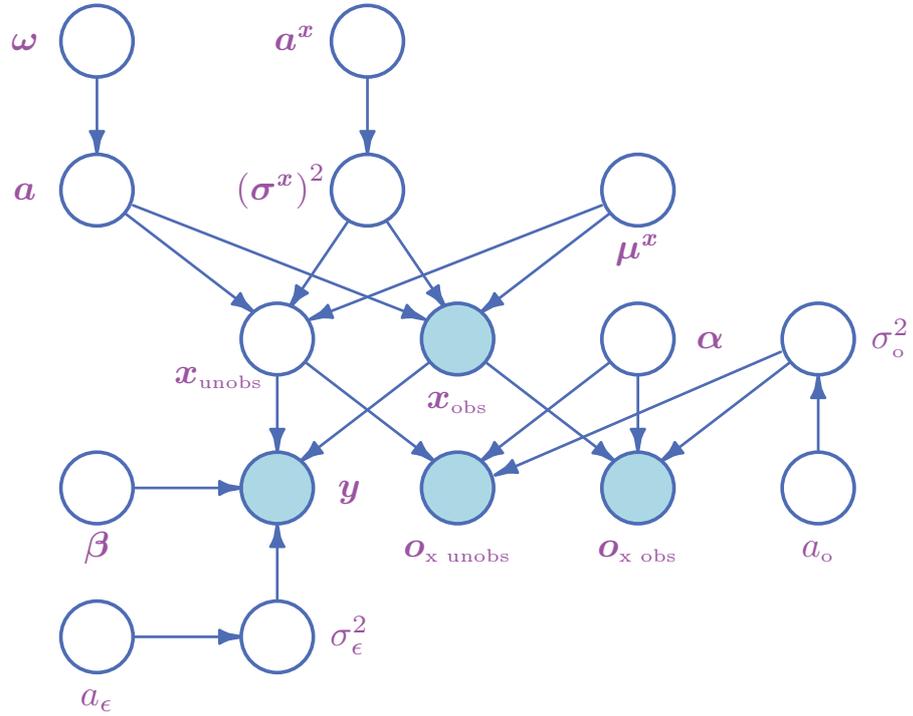


Figure 6.1: DAG corresponding to the model in (6.4). The shaded nodes correspond to observed data and the open nodes correspond to hidden or latent variables.

are defined similarly. It is easier to work with a re-ordered set of variables, that is

$$\mathbf{y} \equiv \begin{bmatrix} \mathbf{y}_{x_{\text{obs}}} \\ \mathbf{y}_{x_{\text{unobs}}} \end{bmatrix}, \quad \mathbf{c} \equiv \begin{bmatrix} \mathbf{c}_{x_{\text{obs}}} \\ \mathbf{c}_{x_{\text{unobs}}} \end{bmatrix}, \quad \text{and} \quad \mathbf{x} \equiv \begin{bmatrix} \mathbf{x}_{\text{obs}} \\ \mathbf{x}_{\text{unobs}} \end{bmatrix}.$$

The conditional dependence structure of (6.4) can be visualised in the DAG given in Figure 6.1. This DAG illustrates the relationship between the regression parameters and the unobserved parameters in (6.4). Lastly, we note that the known covariate vector \mathbf{c} is omitted from the DAG since it does not require inference.

6.4 Mean field variational Bayes

We provide detail on MFVB fitting catered to the Normal response regression model setting with measurement error as described in model (6.4). The assumption behind MFVB is that we impose a product restriction on the joint posterior distribution of our

model parameters. For instance, we may impose the product factorisation

$$\begin{aligned}
 p\{\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{\mu}^{\mathbf{x}}, \mathbf{a}^{\mathbf{x}}, \mathbf{a}, \boldsymbol{\omega}, a_{\varepsilon}, a_o, \sigma_{\varepsilon}^2, \sigma_o^2, (\boldsymbol{\sigma}^{\mathbf{x}})^2 \mid \mathbf{y}, \mathbf{x}_{\text{obs}}, \mathbf{o}\} \\
 \approx q\{\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{\mu}^{\mathbf{x}}, \mathbf{a}^{\mathbf{x}}, \mathbf{a}, \boldsymbol{\omega}, a_{\varepsilon}, a_o, \sigma_{\varepsilon}^2, \sigma_o^2, (\boldsymbol{\sigma}^{\mathbf{x}})^2\} \\
 = q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{\mu}^{\mathbf{x}}, \mathbf{a}^{\mathbf{x}}, \mathbf{a}, \boldsymbol{\omega}, a_{\varepsilon}, a_o) q\{\sigma_{\varepsilon}^2, \sigma_o^2, (\boldsymbol{\sigma}^{\mathbf{x}})^2\}.
 \end{aligned} \tag{6.5}$$

Using the concept of moralisation (see Definition 1.8.9) we can further reduce the product restriction in (6.5) to have the form

$$\begin{aligned}
 q(\mathbf{x}_{\text{unobs}}) q(\boldsymbol{\beta}) q(\boldsymbol{\alpha}) \left\{ \prod_{k=1}^K q(\mu_k^{\mathbf{x}}) q(a_k^{\mathbf{x}}) \right\} \left\{ \prod_{i=1}^n q(a_{i1}, \dots, a_{iK}) \right\} q(\omega_1, \dots, \omega_K) \\
 \times q(a_{\varepsilon}) q(a_o) q(\sigma_{\varepsilon}^2) q(\sigma_o^2) \left[\prod_{k=1}^K q\{(\sigma_k^{\mathbf{x}})^2\} \right].
 \end{aligned} \tag{6.6}$$

In Appendix 6.B, we show that the optimal q -densities for the parameters in (6.6) are:

$q^*(\mathbf{x}_{\text{unobs}})$ is the $N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}, \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \mathbf{I}_{n_{\text{unobs}}})$ density function,

$q^*(\boldsymbol{\beta})$ is the $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})})$ density function,

$q^*(\boldsymbol{\alpha})$ is the $N(\boldsymbol{\mu}_{q(\boldsymbol{\alpha})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\alpha})})$ density function,

$q^*(\boldsymbol{\mu}^{\mathbf{x}}) \equiv \prod_{k=1}^K q(\mu_k^{\mathbf{x}})$ is the product of $N(\mu_{q(\mu_k^{\mathbf{x}})}, \sigma_{q(\mu_k^{\mathbf{x}})}^2)$ density functions,

$q^*(\mathbf{a}^{\mathbf{x}}) \equiv \prod_{k=1}^K q(a_k^{\mathbf{x}})$ is the product of Inverse-Gamma $(1, B_{q(a_k^{\mathbf{x}})})$ density functions,

$q^*(\mathbf{a}) \equiv \prod_{i=1}^n q(a_{i1}, \dots, a_{iK})$ is the product of Multinomial $(1, \mu_{q(a_{i1})}, \dots, \mu_{q(a_{iK})})$ density functions,

$q^*(\omega_1, \dots, \omega_K)$ is the Dirichlet $(\alpha_{q(\omega_1)}, \dots, \alpha_{q(\omega_K)})$ density function,

$q^*(a_{\varepsilon})$ is the Inverse-Gamma $(1, B_{q(a_{\varepsilon})})$ density function,

$q^*(a_o)$ is the Inverse-Gamma $(1, B_{q(a_o)})$ density function,

$q^*(\sigma_{\varepsilon}^2)$ is the Inverse-Gamma $(\frac{1}{2}(n+1), B_{q(\sigma_{\varepsilon}^2)})$ density function,

$q^*(\sigma_o^2)$ is the Inverse-Gamma $(\frac{1}{2}(n+1), B_{q(\sigma_o^2)})$ density function and

$q^*\{(\boldsymbol{\sigma}^{\mathbf{x}})^2\} \equiv q^*\{(\sigma_k^{\mathbf{x}})^2, \dots, (\sigma_k^{\mathbf{x}})^2\}$ is the product of

Inverse-Gamma $(A_{\{(\sigma_k^{\mathbf{x}})^2\}}, B_{\{(\sigma_k^{\mathbf{x}})^2\}})$ density functions.

The optimal parameters in these q -densities possess interdependencies between each other, resulting in an iterative scheme for their solution. This is encompassed in Algorithm 9.

The lower bound on the marginal log-likelihood for model (6.4) is given by (the derivation

of which is given in Appendix 6.C):

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) = & \\
 & -\frac{1}{2}(K+5)\log(\pi) + \frac{1}{2}(p+q+K+n_{\text{unobs}}) - \frac{p}{2}\log(\sigma_\beta^2) + \frac{1}{2}\log|\Sigma_{q(\beta)}| \\
 & -\log(A_\varepsilon) - \frac{1}{2}\left(2n+n_{\text{unobs}} + \sum_{k=1}^K \sum_{i=1}^n \mu_{q(a_{ik})}\right)\log(2\pi) + \frac{1}{2}\log|\Sigma_{q(\alpha)}| \\
 & -\frac{1}{2\sigma_\beta^2}\left\{\|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\Sigma_{q(\beta)})\right\} + \log\Gamma\left\{\frac{1}{2}(n+1)\right\} - \frac{q}{2}\log(\sigma_\alpha^2) \\
 & -\frac{1}{2}(n+1)\log(B_{q(\sigma_\varepsilon^2)}) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/\sigma_\varepsilon^2)}\mu_{q(1/a_\varepsilon)} - \log(A_o) \\
 & -\frac{1}{2\sigma_\alpha^2}\left\{\|\boldsymbol{\mu}_{q(\alpha)}\|^2 + \text{tr}(\Sigma_{q(\alpha)})\right\} - \sum_{k=1}^K \log(A_k^x) + \log\Gamma\left\{\frac{1}{2}(n+1)\right\} \\
 & -\frac{K}{2}\log(\sigma_\mu^2) - \frac{1}{2}(n+1)\log(B_{q(\sigma_o^2)}) - \log(B_{q(a_o)}) + \mu_{q(1/\sigma_o^2)}\mu_{q(1/a_o)} \\
 & + \frac{n_{\text{unobs}}}{2}\log(\sigma_{q(\mathbf{x}_{\text{unobs}})}^2) - \frac{1}{2\sigma_\mu^2}\sum_{k=1}^K\left\{(\mu_{q(\mu_k^x)} - \mu_\mu)^2 + \sigma_{q(\mu_k^x)}^2\right\} + \frac{1}{2}\sum_{k=1}^K \log(\sigma_{q(\mu_k^x)}^2) \\
 & + \sum_{k=1}^K \log\Gamma\left\{\frac{1}{2}(\mu_{q(a_{\bullet k})} + 1)\right\} - \frac{1}{2}\sum_{k=1}^K (\mu_{q(a_{\bullet k})} + 1)\log\left(B_{q\{(\sigma_k^x)^2\}}\right) \\
 & + \sum_{k=1}^K \mu_{q\{1/(\sigma_k^x)^2\}}\mu_{q(1/a_k^x)} - \sum_{i=1}^n \sum_{k=1}^K \log(\mu_{q(a_{ik})})\mu_{q(a_{ik})} - \log\Gamma\left(\sum_{k=1}^K \alpha_{q(\omega_k)}\right) \\
 & + \log\Gamma\left(\sum_{k=1}^K \alpha_k\right) - \sum_{k=1}^K \log\Gamma(\alpha_k) + \sum_{k=1}^K \log\Gamma(\alpha_{q(\omega_k)}) - \sum_{k=1}^K \log(B_{q(a_k^x)}).
 \end{aligned} \tag{6.7}$$

6.5 Simulation study

We conducted a simulation study for the measurement error model given in (6.4), where each i th response measurement takes on a Normal distribution and we make use of a normal mixture model with $K = 2$ mixture components for the prior model on \mathbf{x} . The true parameter values were specified as

$$\begin{aligned}
 \beta_0 = -0.18, \quad \beta_c = 0.44, \quad \beta_x = 0.79, \quad \sigma_\varepsilon^2 = 0.2, \quad \alpha_o = 0.24, \quad \alpha_1 = 0.83, \quad \sigma_o^2 = 0.33, \\
 \omega_1 = 0.25, \quad \omega_2 = 0.75, \quad \mu_1^x = -1.3, \quad \mu_2^x = 1.25, \quad (\sigma_1^x)^2 = 0.45, \quad (\sigma_2^x)^2 = 0.45,
 \end{aligned}$$

and the proportion of observed x 's, p_{obs} , was varied using $p_{\text{obs}} \in \{0.3, 0.5, 0.7, 0.9\}$. We generated 25 datasets for each of these settings. Algorithm 9 was implemented in the R programming language and MCMC was used through the BUGS inference engine. These computations were performed on a laptop computer (Mac OS X; 2.8 GHz processor, 16 GBytes of random access memory). Figure 6.2 shows MFVB accuracy summaries for different types of p_{obs} values. The parameters that are monitored are σ_ε^2 , σ_o^2 and the

Set up initial values:

$$\begin{aligned} & \mu_{q(1/\sigma_o^2)}, \mu_{q(1/a_o)}, \mu_{q(a_{\bullet k})}, \mu_{q(1/(\sigma_k^x)^2)}, \mu_{q(\mu_k^x)}, \sigma_{q(\mu_k^x)}^2, \alpha_{q(\omega_k)}, A_{q((\sigma_k^x)^2)}, A_{q((\mu_k^x)^2)} > 0, \\ & \mu_{q(1/\sigma_\varepsilon^2)}, \mu_{q(1/a_\varepsilon)} > 0. \end{aligned}$$

Cycle through:

Update the expectations given in Appendix 6.A.

Regression model updates:

$$\begin{aligned} \Sigma_{q(\beta)} & \leftarrow \left\{ \mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}^\top \mathbf{X}) + \frac{1}{\sigma_\beta^2} \mathbf{I}_p \right\}^{-1} \\ \boldsymbol{\mu}_{q(\beta)} & \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} \Sigma_{q(\beta)} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X})^\top \mathbf{y} \\ B_{q(\sigma_\varepsilon^2)} & \leftarrow \frac{1}{2} \left(\|\mathbf{y}_{\mathbf{x}_{\text{obs}}} - \mathbf{X}_{\mathbf{x}_{\text{obs}}} \boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\mathbf{X}_{\mathbf{x}_{\text{obs}}}^\top \mathbf{X}_{\mathbf{x}_{\text{obs}}} \Sigma_{q(\beta)}) + \|\mathbf{y}_{\mathbf{x}_{\text{unobs}}}\|^2 \right. \\ & \quad \left. - 2\mathbf{y}_{\mathbf{x}_{\text{unobs}}}^\top E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}_{\mathbf{x}_{\text{unobs}}}) \boldsymbol{\mu}_{q(\beta)} \right. \\ & \quad \left. + \text{tr} \left[E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}_{\mathbf{x}_{\text{unobs}}}^\top \mathbf{X}_{\mathbf{x}_{\text{unobs}}}) \left(\Sigma_{q(\beta)} + \boldsymbol{\mu}_{q(\beta)} \boldsymbol{\mu}_{q(\beta)}^\top \right) \right] \right) + \mu_{q(1/a_\varepsilon)} \\ \mu_{q(1/\sigma_\varepsilon^2)} & \leftarrow \frac{1}{2} (n+1) / B_{q(\sigma_\varepsilon^2)}; \quad B_{q(a_\varepsilon)} \leftarrow \mu_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}; \quad \mu_{q(1/a_\varepsilon)} \leftarrow 1/B_{q(a_\varepsilon)} \end{aligned}$$

Measurement model updates:

$$\begin{aligned} \Sigma_{q(\alpha)} & \leftarrow \left\{ \mu_{q(1/\sigma_o^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) + \frac{1}{\sigma_\alpha^2} \mathbf{I}_q \right\}^{-1} \\ \boldsymbol{\mu}_{q(\alpha)} & \leftarrow \mu_{q(1/\sigma_o^2)} \Sigma_{q(\alpha)} E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}})^\top \mathbf{o} \\ B_{q(\sigma_o^2)} & \leftarrow \frac{1}{2} \left(\|\mathbf{o}_{\mathbf{x}_{\text{obs}}} - \tilde{\mathbf{X}}_{\mathbf{x}_{\text{obs}}} \boldsymbol{\mu}_{q(\alpha)}\|^2 + \text{tr}(\tilde{\mathbf{X}}_{\mathbf{x}_{\text{obs}}}^\top \tilde{\mathbf{X}}_{\mathbf{x}_{\text{obs}}} \Sigma_{q(\alpha)}) + \|\mathbf{o}_{\mathbf{x}_{\text{unobs}}}\|^2 \right. \\ & \quad \left. - 2\mathbf{o}_{\mathbf{x}_{\text{unobs}}}^\top E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}}_{\mathbf{x}_{\text{unobs}}}) \boldsymbol{\mu}_{q(\alpha)} \right. \\ & \quad \left. + \text{tr} \left[E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}}_{\mathbf{x}_{\text{unobs}}}^\top \tilde{\mathbf{X}}_{\mathbf{x}_{\text{unobs}}}) \left(\Sigma_{q(\alpha)} + \boldsymbol{\mu}_{q(\alpha)} \boldsymbol{\mu}_{q(\alpha)}^\top \right) \right] \right) + \mu_{q(1/a_o)} \\ \mu_{q(1/\sigma_o^2)} & \leftarrow \frac{1}{2} (n+1) / B_{q(\sigma_o^2)}; \quad B_{q(a_o)} \leftarrow \mu_{q(1/\sigma_o^2)} + A_o^{-2}; \quad \mu_{q(1/a_o)} \leftarrow 1/B_{q(a_o)} \end{aligned}$$

Prior specification updates:

$$\begin{aligned} \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 & \leftarrow 1 / \left(\mu_{q(1/\sigma_\varepsilon^2)} \left[\left\{ (\boldsymbol{\mu}_{q(\beta)})_3 \right\}^2 + (\Sigma_{q(\beta)})_{33} \right] + \sum_{k=1}^K \mu_{q(a_{\bullet k})} \mu_{q(1/(\sigma_k^x)^2)} \right. \\ & \quad \left. + \mu_{q(1/\sigma_o^2)} \left[\left\{ (\boldsymbol{\mu}_{q(\alpha)})_2 \right\}^2 + (\Sigma_{q(\alpha)})_{22} \right] \right) \end{aligned}$$

$$\begin{aligned}
 \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \leftarrow & \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \left(\mu_{q(1/\sigma_\varepsilon^2)} \left[(\boldsymbol{\mu}_{q(\beta)})_3 \mathbf{y}_{\mathbf{x}_{\text{unobs}}} \right. \right. \\
 & - \mathbf{1}_{n_{\text{unobs}}} \left\{ (\boldsymbol{\mu}_{q(\beta)})_1 + (\boldsymbol{\mu}_{q(\beta)})_3 + (\boldsymbol{\Sigma}_{q(\beta)})_{13} \right\} \\
 & - \mathbf{c}_{\mathbf{x}_{\text{unobs}}} \left\{ (\boldsymbol{\mu}_{q(\beta)})_2 + (\boldsymbol{\mu}_{q(\beta)})_3 + (\boldsymbol{\Sigma}_{q(\beta)})_{23} \right\} \\
 & + \mathbf{1}_{n_{\text{unobs}}} \sum_{k=1}^K \mu_{q(a_{\bullet k})} \mu_{q(\mu_k^x)} \mu_{q(1/(\sigma_k^x)^2)} \\
 & + \mu_{q(1/\sigma_\alpha^2)} \left[(\boldsymbol{\mu}_{q(\alpha)})_1 \mathbf{o}_{\mathbf{x}_{\text{unobs}}} - \mathbf{1}_{n_{\text{unobs}}} \left\{ (\boldsymbol{\mu}_{q(\alpha)})_1 + (\boldsymbol{\mu}_{q(\alpha)})_2 \right. \right. \\
 & \left. \left. + (\boldsymbol{\Sigma}_{q(\alpha)})_{12} \right\} \right] \left. \right],
 \end{aligned}$$

For $i = 1, \dots, n$ and $k = 1, \dots, K$:

$$\begin{aligned}
 \nu_{ik} \leftarrow & \psi(\alpha_{q(\omega_k)}) + \frac{1}{2} \psi \left(A_{q((\sigma_k^x)^2)} \right) - \frac{1}{2} \log \left(B_{q((\sigma_k^x)^2)} \right) \\
 & - \frac{1}{2} \mu_{q(1/(\sigma_k^x)^2)} \left(\left[\{E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x})\}_i - \mu_{q(\mu_k^x)} \right]^2 + \left\{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \right\}_i + \sigma_{q(\mu_k^x)}^2 \right)
 \end{aligned}$$

For $i = 1, \dots, n$ and $k = 1, \dots, K$: $\mu_{q(a_{ik})} \leftarrow \exp(\nu_{ik}) / \sum_{k=1}^K \exp(\nu_{ik})$

For $k = 1, \dots, K$:

$$\begin{aligned}
 \mu_{q(a_{\bullet k})} \leftarrow & \sum_{i=1}^n \mu_{q(a_{ik})}; \quad \sigma_{q(\mu_k^x)}^2 \leftarrow 1 / \left\{ \mu_{q(1/(\sigma_k^x)^2)} \mu_{q(a_{\bullet k})} + 1/\sigma_\mu^2 \right\} \\
 \mu_{q(\mu_k^x)} \leftarrow & \sigma_{q(\mu_k^x)}^2 \left\{ \mu_{q(1/(\sigma_k^x)^2)} \mu_{q(a_{\bullet k})} \sum_{i=1}^n \left\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}) \right\}_i + \mu_\mu / \sigma_\mu^2 \right\} \\
 \alpha_{q(\omega_k)} \leftarrow & \mu_{q(a_{\bullet k})} + \alpha; \quad A_{q((\sigma_k^x)^2)} \leftarrow \frac{1}{2} (\mu_{q(a_{\bullet k})} + 1) \\
 B_{q((\sigma_k^x)^2)} \leftarrow & \frac{1}{2} \mu_{q(a_{\bullet k})} \sum_{i=1}^n \left(\left[\{E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x})\}_i - \mu_{q(\mu_k^x)} \right]^2 + \left\{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \right\}_i + \sigma_{q(\mu_k^x)}^2 \right) \\
 & + \mu_{q(1/a_k^x)} \\
 \mu_{q(1/(\sigma_k^x)^2)} \leftarrow & A_{q((\sigma_k^x)^2)} / B_{q((\sigma_k^x)^2)}; \quad B_{q(a_k^x)} \leftarrow \mu_{q(1/(\sigma_k^x)^2)} + (A_k^x)^{-2} \\
 \mu_{q(1/a_k^x)} \leftarrow & 1 / B_{q(a_k^x)}
 \end{aligned}$$

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

Algorithm 9: *MFVB algorithm for obtaining the parameters in the optimal q densities for the normal response measurement error model (6.4).*

6.5. SIMULATION STUDY

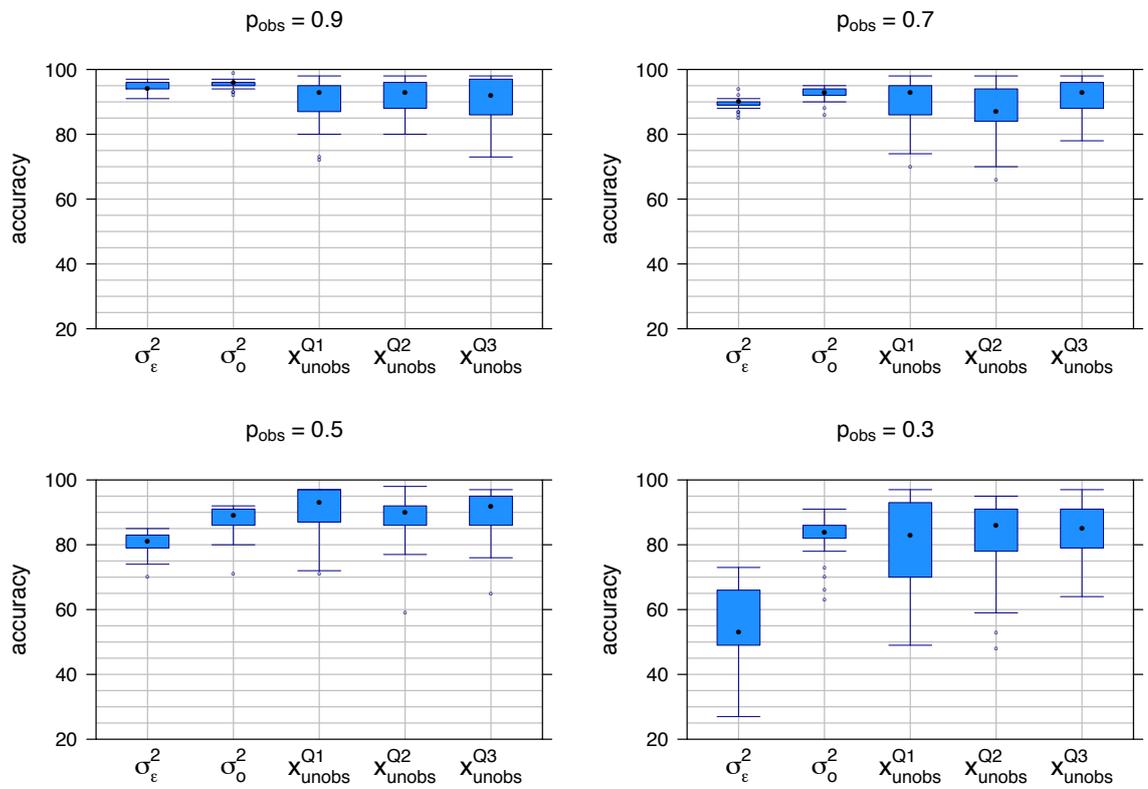


Figure 6.2: *Boxplots of parameter accuracies for MFVB in Algorithm 9 applied to the measurement error model given in (6.4) for the simulation setting described in Section 6.5.*

6.5. SIMULATION STUDY

quartiles of the unobserved x 's. These are represented as $x_{\text{unobs}}^{\text{Q}_i}$, $i = 1, 2, 3$. We can see from Figure 6.2 that the accuracies seem to be quite good for σ_o^2 and the quartiles of $\mathbf{x}_{\text{unobs}}$ for all values of p_{obs} except for $p_{\text{obs}} = 0.3$. However, accuracies are very poor for σ_ε^2 when $p_{\text{obs}} = 0.3$. Overall, the accuracies are looking good when $p_{\text{obs}} \in \{0.5, 0.7, 0.9\}$.

We also assessed the 95 % credible interval coverage achieved by the MFVB optimal q -densities against the true parameter coverage for the same parameters. From Table 6.1 we see that the coverage probabilities are generally high for all parameters except σ_ε^2 and σ_o^2 . As p_{obs} decreases, so do the coverage probabilities for σ_ε^2 and σ_o^2 . The run time for

p_{obs}	0.9	0.7	0.5	0.3
σ_ε^2	92	96	76	68
σ_o^2	92	88	84	72
$x_{\text{unobs}}^{\text{Q}_1}$	96	84	96	84
$x_{\text{unobs}}^{\text{Q}_2}$	96	92	100	100
$x_{\text{unobs}}^{\text{Q}_3}$	92	96	100	92

Table 6.1: 95% credible interval coverage probabilities for the MFVB approximation applied to (6.4).

p_{obs}	0.9	0.7	0.5	0.3
MCMC	51.80 (0.61)	55.69 (0.51)	54.46 (1.70)	59.80 (0.67)
MFVB	1.46 (0.10)	1.46 (0.11)	1.42 (0.03)	1.41 (0.04)
Ratio	35.48	38.17	38.44	42.46

Table 6.2: Average (standard deviation) run time in seconds for MFVB and MCMC fitting of the Gaussian response measurement error model.

MCMC and MFVB was also assessed. The resulting times are given in table 6.2. Based on the timing results, we see that MFVB is approximately 40 times faster than MCMC. The time savings afforded by MFVB in this instance is not of a large magnitude compared with other models considered in this thesis, however extension to arbitrarily large and complex models involving measurement error would more clearly represent the time benefits of using MFVB as an alternative to MCMC.

From Figure 6.3 we can see that the means of the MFVB fits are agreeing well to the fits obtained via MCMC. More importantly, both of these fits are agreeing with the true mean function even when $p_{\text{obs}} = 0.3$. Thus, overall this simulation study has represented the effectiveness of using MFVB as an alternative to MCMC in models impaired

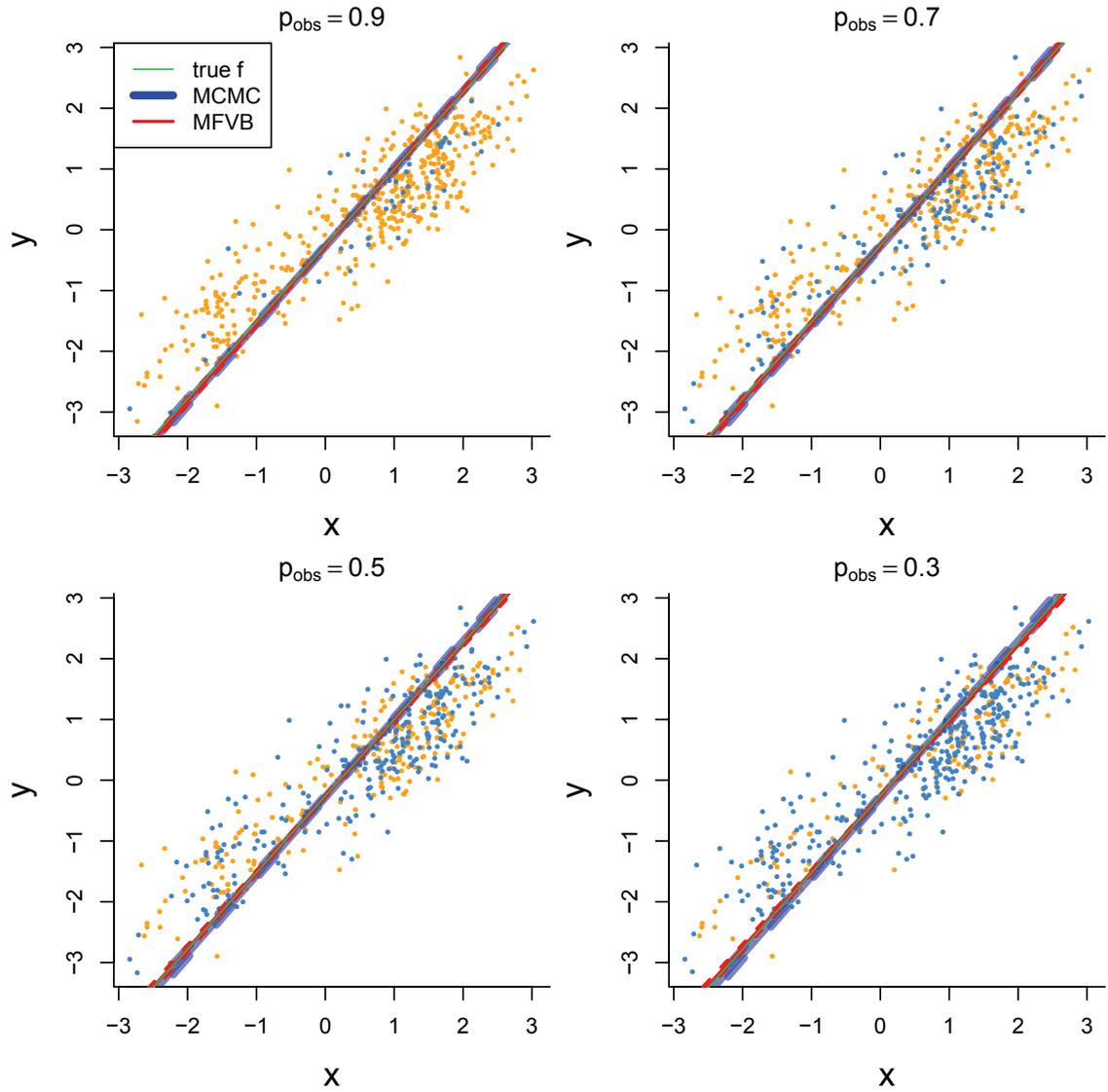


Figure 6.3: Comparative plots of the MFVB and MCMC mean fitted functions with corresponding 95% credible sets for a single replication of the simulation study. The orange data points represent the observed data and the blue data points represent the unobserved data.

by measurement error.

6.6 Discussion

We have derived a fast MFVB algorithm to a regression model with classical measurement error. This approach was shown to have good accuracy with evidence of time saving compared to the MCMC approach. Even though the time comparisons weren't of high magnitude for MFVB compared to MCMC, extension of such methodology to arbitrarily large and complex models involving measurement error would exhibit larger time savings.

6.A Expectation updates

This section lists the updates in Algorithm 9 that involve the expectation operator with respect to $\mathbf{x}_{\text{unobs}}$:

$$\begin{aligned}
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}) &\leftarrow \begin{bmatrix} \mathbf{x}_{\text{obs}} \\ \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}^\top \mathbf{x}) &\leftarrow \|\mathbf{x}_{\text{obs}}\|^2 + \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}})\|^2 + \sum_{i=1}^{n_{\text{unobs}}} \sigma_q^2(\mathbf{x}_{\text{unobs}}), \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\tilde{\mathbf{X}}) &\leftarrow \begin{bmatrix} \mathbf{1} & \mathbf{x}_{\text{obs}} \\ \mathbf{1} & \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) &\leftarrow \begin{bmatrix} n & \mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) & E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}^\top \mathbf{x}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\tilde{\mathbf{X}} \mathbf{x}_{\text{unobs}}) &\leftarrow \begin{bmatrix} \mathbf{1} & \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\tilde{\mathbf{X}}^\top \mathbf{x}_{\text{unobs}} \tilde{\mathbf{X}} \mathbf{x}_{\text{unobs}}) &\leftarrow \begin{bmatrix} n_{\text{unobs}} & \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) & \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}})\|^2 + \sum_{i=1}^{n_{\text{unobs}}} \sigma_q^2(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{X}) &\leftarrow \begin{bmatrix} \mathbf{1} & \mathbf{c}_{\text{obs}} & \mathbf{x}_{\text{obs}} \\ \mathbf{1} & \mathbf{c}_{\text{unobs}} & \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{X}^\top \mathbf{X}) &\leftarrow \begin{bmatrix} n & \mathbf{1}^\top \mathbf{c} \\ \mathbf{1}^\top \mathbf{c} & \mathbf{c}^\top \mathbf{c} \\ \mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) & \mathbf{c}_{\text{obs}}^\top \mathbf{x}_{\text{obs}} + \mathbf{c}_{\text{unobs}}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) & \mathbf{1}^\top \mathbf{x}_{\text{obs}} + \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ & & \mathbf{c}_{\text{obs}}^\top \mathbf{x}_{\text{obs}} + \mathbf{c}_{\text{unobs}}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ & & E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}^\top \mathbf{x}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{X} \mathbf{x}_{\text{unobs}}) &\leftarrow \begin{bmatrix} \mathbf{1} & \mathbf{c}_{\text{unobs}} & \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \end{bmatrix}, \\
 E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{X}_{\text{unobs}}^\top \mathbf{X}_{\text{unobs}}) &\leftarrow \begin{bmatrix} n_{\text{unobs}} & \mathbf{1}^\top \mathbf{c}_{\text{unobs}} \\ \mathbf{1}^\top \mathbf{c}_{\text{unobs}} & \|\mathbf{c}_{\text{unobs}}\|^2 \\ \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) & \mathbf{c}_{\text{unobs}}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ & & \mathbf{1}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ & & \mathbf{c}_{\text{unobs}}^\top \boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}}) \\ & & \|\boldsymbol{\mu}_q(\mathbf{x}_{\text{unobs}})\|^2 + \sum_{i=1}^{n_{\text{unobs}}} \sigma_q^2(\mathbf{x}_{\text{unobs}}) \end{bmatrix}.
 \end{aligned}$$

Also,

$$\boldsymbol{\sigma}_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \equiv \begin{bmatrix} \mathbf{0}_{n_{\text{obs}}} \\ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \mathbf{1}_{n_{\text{unobs}}} \end{bmatrix},$$

where $\mathbf{0}_{n_{\text{obs}}}$ is the $n_{\text{obs}} \times 1$ column vector with all entries equal to 0.

6.B Derivation of Algorithm 9

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

First note that

$$\begin{aligned} p(\boldsymbol{\beta} \mid \text{rest}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) p(\boldsymbol{\beta}) \\ &= \exp \left\{ -\frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{1}{2\sigma_\beta^2} \|\boldsymbol{\beta}\|^2 \right\} + \text{const}, \end{aligned}$$

where ‘const’ denotes all terms not depending on $\boldsymbol{\beta}$. Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\boldsymbol{\beta} \mid \text{rest}) &\propto -\frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{1}{2\sigma_\beta^2} \|\boldsymbol{\beta}\|^2 \\ &\propto -\frac{1}{2\sigma_\varepsilon^2} (\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}) - \frac{1}{2\sigma_\beta^2} \|\boldsymbol{\beta}\|^2 \\ &= -\frac{1}{2} \left\{ \boldsymbol{\beta}^\top \left(\frac{1}{\sigma_\varepsilon^2} \mathbf{X}^\top \mathbf{X} + \frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right) \boldsymbol{\beta} - 2\boldsymbol{\beta}^\top \left(\frac{1}{\sigma_\varepsilon^2} \mathbf{X}^\top \mathbf{y} \right) \right\} + \text{const}. \end{aligned}$$

Taking expectations with respect to all parameters except $\boldsymbol{\beta}$, we get

$$\begin{aligned} \log q^*(\boldsymbol{\beta}) &= E_q \{ \log p(\boldsymbol{\beta} \mid \text{rest}) \} + \text{const}. \\ &= -\frac{1}{2} \left\{ \boldsymbol{\beta}^\top \left(\mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}^\top \mathbf{X}) + \frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right) \boldsymbol{\beta} \right. \\ &\quad \left. - 2\boldsymbol{\beta}^\top \left(\mu_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}^\top \mathbf{y}) \right) \right\} + \text{const}. \end{aligned}$$

Therefore,

$$q^*(\boldsymbol{\beta}) \text{ is the } N(\boldsymbol{\mu}_{q(\boldsymbol{\beta})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \text{ density function,}$$

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\boldsymbol{\beta})} &= \boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}^\top \mathbf{y}), \text{ and} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} &= \left(\boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\mathbf{X}^\top \mathbf{X}) + \frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right)^{-1}. \end{aligned}$$

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

$$\begin{aligned}
 p(\sigma_\varepsilon^2 | \text{rest}) &\propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) p(\sigma_\varepsilon^2 | a_\varepsilon) \\
 &= |\sigma_\varepsilon^2 \mathbf{I}_n|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right\} (\sigma_\varepsilon^2)^{-\frac{1}{2}-1} \exp\left\{-(1/a_\varepsilon)/\sigma_\varepsilon^2\right\} + \text{const.}
 \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned}
 \log p(\sigma_\varepsilon^2 | \text{rest}) &= -\frac{n}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{3}{2} \log(\sigma_\varepsilon^2) - (1/a_\varepsilon)/\sigma_\varepsilon^2 \\
 &\quad + \text{const.} \\
 &= -\left\{\frac{1}{2}(n+1) + 1\right\} \log(\sigma_\varepsilon^2) - \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + 1/a_\varepsilon\right) / \sigma_\varepsilon^2 + \text{const.}
 \end{aligned}$$

Taking expectations with respect to all parameters except σ_ε^2 , we get

$$\begin{aligned}
 \log q^*(\sigma_\varepsilon^2) &= E_q \left\{ \log q^*(\sigma_\varepsilon^2 | \text{rest}) \right\} + \text{const.} \\
 &= -\left\{\frac{1}{2}(n+1) + 1\right\} \log(\sigma_\varepsilon^2) - \left(\frac{1}{2} E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \boldsymbol{\mu}_{q(1/a_\varepsilon)}\right) / \sigma_\varepsilon^2 \\
 &\quad + \text{const.}
 \end{aligned}$$

Now, breaking $E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ up into the *observed* and *unobserved* parts, we get

$$E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = E_q \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{X}_{x_{\text{obs}}}\boldsymbol{\beta}\|^2 + E_q \|\mathbf{y}_{x_{\text{unobs}}} - \mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta}\|^2.$$

Using Result 1.4.19,

$$\begin{aligned}
 E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 &= \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{X}_{x_{\text{obs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}\left(\mathbf{X}_{x_{\text{obs}}}^\top \mathbf{X}_{x_{\text{obs}}}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}\right) \\
 &\quad + \|\mathbf{y}_{x_{\text{unobs}}} - E_q(\mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}\{\text{Cov}_q(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta})\}
 \end{aligned}$$

where, using Result 1.4.20,

$$\begin{aligned}
 \text{Cov}_q(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta}) &= E_q \left\{ \text{Cov}_q(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta} | \mathbf{x}_{\text{unobs}}) \right\} + \text{Cov}_q \left\{ E_q(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta} | \mathbf{x}_{\text{unobs}}) \right\} \\
 &= E_q \left\{ \mathbf{X}_{x_{\text{unobs}}}\text{Cov}_q(\boldsymbol{\beta} | \mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}^\top \right\} + \text{Cov}_q \left\{ \mathbf{X}_{x_{\text{unobs}}} E_q(\boldsymbol{\beta} | \mathbf{x}_{\text{unobs}}) \right\}
 \end{aligned}$$

Taking into account the independence between $\mathbf{x}_{\text{unobs}}$ and $\boldsymbol{\beta}$, as shown in product restriction (6.6), this gives

$$\begin{aligned}
 \text{Cov}_q(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta}) &= E_q \left\{ \mathbf{X}_{x_{\text{unobs}}}\text{Cov}_q(\boldsymbol{\beta})\mathbf{X}_{x_{\text{unobs}}}^\top \right\} + \text{Cov}_q \left\{ \mathbf{X}_{x_{\text{unobs}}} E_q(\boldsymbol{\beta}) \right\} \\
 &= E_q \left(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}\mathbf{X}_{x_{\text{unobs}}}^\top \right) + \text{Cov}_q \left(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})} \right) \\
 &= E_q \left(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}\mathbf{X}_{x_{\text{unobs}}}^\top \right) + E_q \left(\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\boldsymbol{\mu}_{q(\boldsymbol{\beta})}^\top \mathbf{X}_{x_{\text{unobs}}}^\top \right) \\
 &\quad - \left\{ E_q(\mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})} \right\} \left\{ E_q(\mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})} \right\}^\top \\
 &= E_q(\mathbf{x}_{\text{unobs}}) \left(\mathbf{X}_{x_{\text{unobs}}}^\top \mathbf{X}_{x_{\text{unobs}}} \right) \left(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta})}\boldsymbol{\mu}_{q(\boldsymbol{\beta})}^\top \right) \\
 &\quad - \left\{ E_q(\mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})} \right\} \left\{ E_q(\mathbf{x}_{\text{unobs}})\mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\mu}_{q(\boldsymbol{\beta})} \right\}^\top.
 \end{aligned}$$

Therefore,

$$\begin{aligned} E_q \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|^2 &= \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{X}_{x_{\text{obs}}} \boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\mathbf{X}_{x_{\text{obs}}}^\top \mathbf{X}_{x_{\text{obs}}} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \\ &\quad + \|\mathbf{y}_{x_{\text{unobs}}}\|^2 - 2\mathbf{y}_{x_{\text{unobs}}}^\top E_q(\mathbf{x}_{\text{unobs}}) (\mathbf{X}_{x_{\text{unobs}}}) \boldsymbol{\mu}_{q(\boldsymbol{\beta})} \\ &\quad + \text{tr} \left[E_q(\mathbf{x}_{\text{unobs}}) (\mathbf{X}_{x_{\text{unobs}}}^\top \mathbf{X}_{x_{\text{unobs}}}) (\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})} + \boldsymbol{\mu}_{q(\boldsymbol{\beta})} \boldsymbol{\mu}_{q(\boldsymbol{\beta})}^\top) \right]. \end{aligned}$$

And so,

$$q^*(\sigma_\varepsilon^2) \text{ follows an Inverse-Gamma } \left(\frac{1}{2}(n+1), B_{q(\sigma_\varepsilon^2)}\right) \text{ distribution,}$$

where

$$B_{q(\sigma_\varepsilon^2)} = \frac{1}{2} E_q \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|^2 + \boldsymbol{\mu}_{q(1/a_\varepsilon)}.$$

Expressions for $B_{q(a_\varepsilon)}$ and $\boldsymbol{\mu}_{q(1/a_\varepsilon)}$

We begin with

$$\begin{aligned} p(a_\varepsilon | \text{rest}) &\propto p(\sigma_\varepsilon^2 | a_\varepsilon) p(a_\varepsilon) \\ &= (1/a_\varepsilon)^{1/2} \exp\{-(1/a_\varepsilon)/\sigma_\varepsilon^2\} (a_\varepsilon)^{-\frac{1}{2}-1} \exp\{-(1/A_\varepsilon^2/a_\varepsilon)\} + \text{const.} \end{aligned}$$

Taking the logarithm of both sides, gives

$$\begin{aligned} \log p(a_\varepsilon | \text{rest}) &= -\frac{1}{2} \log(a_\varepsilon) - (1/a_\varepsilon)/\sigma_\varepsilon^2 - \frac{3}{2} \log(a_\varepsilon) - (1/A_\varepsilon^2)/a_\varepsilon + \text{const} \\ &= -2 \log(a_\varepsilon) - (\sigma_\varepsilon^{-2} + A_\varepsilon^{-2})/a_\varepsilon + \text{const.} \end{aligned}$$

Taking expectations, we get

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

Therefore,

$$q^*(a_\varepsilon) \text{ is the Inverse-Gamma } (1, B_{q(a_\varepsilon)}) \text{ density function,}$$

where

$$B_{q(a_\varepsilon)} = \boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)} + A_\varepsilon^{-2}.$$

In addition, using Result 1.4.3,

$$\boldsymbol{\mu}_{q(1/a_\varepsilon)} = 1/B_{q(a_\varepsilon)}.$$

Expressions for $\mu_q(\mathbf{x}_{\text{unobs}})$ and $\sigma_q^2(\mathbf{x}_{\text{unobs}})$

We first note that

$$\begin{aligned}
 p(\mathbf{x}_{\text{unobs}} \mid \text{rest}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) p(\mathbf{x} \mid \mathbf{a}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{o} \mid \mathbf{x}, \boldsymbol{\alpha}, \sigma_o^2) \\
 &= \exp\left\{-\frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y}_{x_{\text{unobs}}} - \mathbf{X}_{x_{\text{unobs}}}\boldsymbol{\beta}\|^2\right\} \\
 &\quad \times \prod_{i=n_{\text{obs}}+1}^{n_{\text{unobs}}} \prod_{k=1}^K \left[\{2\pi(\sigma_k^x)^2\}^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \mu_k^x)^2 / (\sigma_k^x)^2\right\} \right]^{a_{ik}} \\
 &\quad \times \exp\left\{-\frac{1}{2\sigma_o^2} \|\mathbf{o}_{x_{\text{unobs}}} - \tilde{\mathbf{X}}_{x_{\text{unobs}}}\boldsymbol{\alpha}\|^2\right\} + \text{const.}
 \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned}
 \log p(\mathbf{x}_{\text{unobs}} \mid \text{rest}) &= -\frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y}_{x_{\text{unobs}}} - \beta_0 - \beta_c \mathbf{c}_{x_{\text{unobs}}} - \beta_x \mathbf{x}_{\text{unobs}}\|^2 \\
 &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{unobs}}} a_{ik} \frac{(\xi_i - \mu_k)^2}{\sigma_k^2} \\
 &\quad - \frac{1}{2\sigma_o^2} \|\mathbf{o}_{x_{\text{unobs}}} - \alpha_0 - \alpha_1 \mathbf{x}_{\text{unobs}}\|^2 + \text{const} \\
 &= -\frac{1}{2\sigma_\varepsilon^2} \left(\beta_x^2 \|\mathbf{x}_{\text{unobs}}\|^2 + 2\beta_0 \beta_x \mathbf{1}_{n_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} - 2\beta_x \mathbf{y}_{x_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} \right. \\
 &\quad \left. + 2\beta_c \beta_x \mathbf{c}_{x_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} \right) - \frac{1}{2} \left(\|\mathbf{x}_{\text{unobs}}\|^2 \sum_{k=1}^K \frac{a_{\cdot k}}{\sigma_k^2} \right. \\
 &\quad \left. - 2 \left(\mathbf{1}_{n_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} \right) \sum_{k=1}^K \frac{a_{\cdot k} \mu_k}{\sigma_k^2} \right) - \frac{1}{2\sigma_o^2} \left(\alpha_1^2 \|\mathbf{x}_{\text{unobs}}\|^2 \right. \\
 &\quad \left. + 2\alpha_0 \alpha_1 \mathbf{1}_{n_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} - 2\alpha_1 \mathbf{o}_{x_{\text{unobs}}}^\top \mathbf{x}_{\text{unobs}} \right) + \text{const} \\
 &= -\frac{1}{2} \mathbf{x}_{\text{unobs}}^\top \left(\frac{\beta_x^2}{\sigma_\varepsilon^2} + \sum_{k=1}^K \frac{a_{\cdot k}}{\sigma_k^2} + \frac{\alpha_1^2}{\sigma_o^2} \right) \mathbf{x}_{\text{unobs}} \\
 &\quad + \mathbf{x}_{\text{unobs}}^\top \left(\frac{\beta_x}{\sigma_\varepsilon^2} (\mathbf{y}_{x_{\text{unobs}}} - \beta_0 \mathbf{1}_{n_{\text{unobs}}} - \beta_c \mathbf{c}_{x_{\text{unobs}}}) \right. \\
 &\quad \left. + \mathbf{1}_{n_{\text{unobs}}} \sum_{k=1}^K \frac{a_{\cdot k} \mu_k}{\sigma_k^2} + \frac{\alpha_1}{\sigma_o^2} (\mathbf{o}_{x_{\text{unobs}}} - \alpha_0 \mathbf{1}_{n_{\text{unobs}}}) \right) + \text{const.}
 \end{aligned}$$

Taking expectations, we get

$$\begin{aligned}
 \log q^*(\mathbf{x}_{\text{unobs}}) &= E_q \{ \log p(\mathbf{x}_{\text{unobs}} \mid \text{rest}) \} + \text{const} \\
 &= -\frac{1}{2} \left[\mathbf{x}_{\text{unobs}}^\top E_q \left(\frac{\beta_x^2}{\sigma_\varepsilon^2} + \sum_{k=1}^K \frac{a_{\cdot k}}{\sigma_k^2} + \frac{\alpha_1^2}{\sigma_o^2} \right) \mathbf{x}_{\text{unobs}} \right. \\
 &\quad \left. - 2 \mathbf{x}_{\text{unobs}}^\top E_q \left(\frac{\beta_x}{\sigma_\varepsilon^2} (\mathbf{y}_{x_{\text{unobs}}} - \beta_0 \mathbf{1}_{n_{\text{unobs}}} - \beta_c \mathbf{c}_{x_{\text{unobs}}}) \right. \right. \\
 &\quad \left. \left. + \mathbf{1}_{n_{\text{unobs}}} \sum_{k=1}^K \frac{a_{\cdot k} \mu_k}{\sigma_k^2} + \frac{\alpha_1}{\sigma_o^2} (\mathbf{o}_{x_{\text{unobs}}} - \alpha_0 \mathbf{1}_{n_{\text{unobs}}}) \right) \right] \\
 &\quad + \text{const.}
 \end{aligned}$$

Therefore,

$q^*(\mathbf{x}_{\text{unobs}})$ is the $N(\boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})}, \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \mathbf{I}_{n_{\text{unobs}}})$ density function,

where

$$\begin{aligned} \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 &= 1/E_q \left(\frac{\beta_x^2}{\sigma_\varepsilon^2} + \sum_{k=1}^K \frac{a_{.k}}{\sigma_k^2} + \frac{\alpha_1^2}{\sigma_o^2} \right) \\ &= 1/\left[\boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)} \left\{ \boldsymbol{\mu}_{q(\beta_x)}^2 + (\boldsymbol{\Sigma}_{q(\beta)})_{33} \right\} + \sum_{k=1}^K \boldsymbol{\mu}_{q(a_{.k})} \boldsymbol{\mu}_{q(1/\sigma_k^2)} \right. \\ &\quad \left. + \boldsymbol{\mu}_{q(1/\sigma_o^2)} \left\{ \boldsymbol{\mu}_{q(\alpha_1)}^2 + (\boldsymbol{\Sigma}_{q(\alpha)})_{22} \right\} \right], \text{ and} \end{aligned}$$

$$\begin{aligned} \boldsymbol{\mu}_{q(\mathbf{x}_{\text{unobs}})} &= E_q \left(\frac{\beta_x}{\sigma_\varepsilon^2} (\mathbf{y}_{x_{\text{unobs}}} - \beta_0 \mathbf{1}_{n_{\text{unobs}}} - \beta_c \mathbf{c}_{x_{\text{unobs}}}) \right. \\ &\quad \left. + \mathbf{1}_{n_{\text{unobs}}} \sum_{k=1}^K \frac{a_{.k} \mu_k}{\sigma_k^2} + \frac{\alpha_1}{\sigma_o^2} (\mathbf{o}_{x_{\text{unobs}}} - \alpha_0 \mathbf{1}_{n_{\text{unobs}}}) \right) \\ &= \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \left[\boldsymbol{\mu}_{q(1/\sigma_\varepsilon^2)} \left[\boldsymbol{\mu}_{q(\beta_x)} \mathbf{y}_{x_{\text{unobs}}} - \mathbf{1}_{n_{\text{unobs}}} (\boldsymbol{\mu}_{q(\beta_0)} \boldsymbol{\mu}_{q(\beta_x)} + (\boldsymbol{\Sigma}_{q(\beta)})_{13}) \right. \right. \\ &\quad \left. \left. - \mathbf{c}_{x_{\text{unobs}}} \left\{ \boldsymbol{\mu}_{q(\beta_c)} \boldsymbol{\mu}_{q(\beta_x)} + (\boldsymbol{\Sigma}_{q(\beta)})_{23} \right\} \right] + \boldsymbol{\mu}_{q(1/\sigma_o^2)} \left\{ \boldsymbol{\mu}_{q(\alpha_1)} \mathbf{o}_{x_{\text{unobs}}} \right. \right. \\ &\quad \left. \left. - \mathbf{1}_{n_{\text{unobs}}} (\boldsymbol{\mu}_{q(\alpha_0)} \boldsymbol{\mu}_{q(\alpha_1)} + (\boldsymbol{\Sigma}_{q(\alpha)})_{12}) \right\} + \mathbf{1}_{n_{\text{unobs}}} \sum_{k=1}^K \boldsymbol{\mu}_{q(a_{.k})} \boldsymbol{\mu}_{q(\mu_k)} \boldsymbol{\mu}_{q(1/\sigma_k^2)} \right]. \end{aligned}$$

Expressions for $\boldsymbol{\mu}_{q(\alpha)}$ and $\boldsymbol{\Sigma}_{q(\alpha)}$

This derivation is similar to the derivation for $\boldsymbol{\mu}_{q(\beta)}$ and $\boldsymbol{\Sigma}_{q(\beta)}$, therefore

$q^*(\alpha)$ is the $N(\boldsymbol{\mu}_{q(\alpha)}, \boldsymbol{\Sigma}_{q(\alpha)})$ density function,

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\alpha)} &= \boldsymbol{\mu}_{q(1/\sigma_o^2)} \boldsymbol{\Sigma}_{q(\alpha)} E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}})^\top \mathbf{o}, \text{ and} \\ \boldsymbol{\Sigma}_{q(\alpha)} &= \left(\boldsymbol{\mu}_{q(1/\sigma_o^2)} E_{q(\mathbf{x}_{\text{unobs}})} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) + \frac{1}{\sigma_\alpha^2} \mathbf{I}_2 \right)^{-1}. \end{aligned}$$

Expressions for $B_{q(\sigma_o^2)}$, $\mu_{q(1/\sigma_o^2)}$, $B_{q(a_o)}$ and $\mu_{q(1/a_o)}$

These derivations are similar to those given for $B_{q(\sigma_\varepsilon^2)}$, $\mu_{q(1/\sigma_\varepsilon^2)}$, $B_{q(a_\varepsilon)}$ and $\mu_{q(1/a_\varepsilon)}$ and so we get the following updates:

$q^*(\sigma_o^2)$ is the Inverse-Gamma $(\frac{1}{2}(n+1), B_{q(\sigma_o^2)})$ density function,

where

$$\begin{aligned}
 B_{q(\sigma_o^2)} &= \frac{1}{2} E_q \left\{ \left\| \mathbf{o} - \tilde{\mathbf{X}} \boldsymbol{\alpha} \right\|^2 \right\} + \mu_{q(1/a_o)} \\
 &= \frac{1}{2} \left[\left\| \mathbf{o}_{x_{\text{obs}}} - \tilde{\mathbf{X}}_{x_{\text{obs}}} \boldsymbol{\mu}_q(\boldsymbol{\alpha}) \right\|^2 + \text{tr} \left(\tilde{\mathbf{X}}_{x_{\text{obs}}} \boldsymbol{\Sigma}_q(\boldsymbol{\alpha}) \tilde{\mathbf{X}}_{x_{\text{obs}}}^\top \right) \right. \\
 &\quad \left. + \left\| \mathbf{o}_{x_{\text{unobs}}} \right\|^2 - 2 \mathbf{o}_{x_{\text{unobs}}} E_{q(\mathbf{x}_{\text{unobs}})} \left(\tilde{\mathbf{X}}_{x_{\text{unobs}}} \right) \boldsymbol{\mu}_q(\boldsymbol{\alpha}) \right. \\
 &\quad \left. + \text{tr} \left\{ E_{q(\mathbf{x}_{\text{unobs}})} \left(\tilde{\mathbf{X}}_{x_{\text{unobs}}}^\top \tilde{\mathbf{X}}_{x_{\text{unobs}}} \right) \left(\boldsymbol{\Sigma}_q(\boldsymbol{\alpha}) + \boldsymbol{\mu}_q(\boldsymbol{\alpha}) \boldsymbol{\mu}_q(\boldsymbol{\alpha})^\top \right) \right\} \right] + \mu_{q(1/\sigma_o^2)}.
 \end{aligned}$$

In addition, using Result 1.4.3, we get

$$\mu_{q(1/\sigma_o^2)} = \frac{1}{2} (n+1) / B_{q(\sigma_o^2)}.$$

Also,

$$q^*(a_o) \text{ is the Inverse-Gamma}(1, B_{q(a_o)}) \text{ density function,}$$

where

$$B_{q(a_o)} = \mu_{q(1/\sigma_o^2)} + A_o^{-2}, \text{ and}$$

$$\mu_{q(1/a_o)} = 1/B_{q(a_o)}.$$

Expressions for $\mu_{q(\mu_k^x)}$ and $\sigma_{q(\mu_k^x)}^2$, $1 \leq k \leq K$

Keeping in mind that $\boldsymbol{\mu}^x = (\mu_1^x, \dots, \mu_K^x)$, we begin with

$$\begin{aligned}
 p(\boldsymbol{\mu}^x | \text{rest}) &\propto p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} p(\boldsymbol{\mu}^x) \\
 &= \prod_{k=1}^K \prod_{i=1}^n \left[\left\{ 2\pi (\sigma_k^x)^2 \right\}^{-1/2} \exp \left\{ -\frac{(x_i - \mu_k^x)^2}{2(\sigma_k^x)^2} \right\} \right]^{a_{ik}} \\
 &\quad \times \prod_{k=1}^K \left[\left(2\pi \sigma_\mu^2 \right)^{-1/2} \exp \left\{ -\frac{(\mu_k^x - \mu_\mu)^2}{\sigma_\mu^2} \right\} \right].
 \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned}
 \log p(\boldsymbol{\mu}^x | \text{rest}) &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n a_{ik} \left\{ \frac{(\mu_k^x)^2 - 2x_i \mu_k^x}{(\sigma_k^x)^2} \right\} - \frac{1}{2\sigma_\mu^2} \left\{ (\mu_k^x)^2 - 2\mu_k^x \mu_\mu \right\} + \text{const} \\
 &= -\frac{1}{2} \sum_{k=1}^K \left[(\mu_k^x)^2 \left\{ \frac{a_{\bullet k}}{(\sigma_k^x)^2} + \frac{1}{\sigma_\mu^2} \right\} - 2\mu_k^x \left\{ \sum_{i=1}^n \frac{\sum_{k=1}^n a_{ik} x_i}{(\sigma_k^x)^2} + \frac{\mu_\mu}{\sigma_\mu^2} \right\} \right].
 \end{aligned}$$

Taking expectations, we get

$$E_q \{ \log(\boldsymbol{\mu} | \text{rest}) \} = -\frac{1}{2} \sum_{k=1}^K \left((\mu_k^x)^2 \left[\mu_{q\{1/(\sigma_k^x)^2\}} \mu_{q(a_{\bullet k})} + \frac{1}{\sigma_\mu^2} \right] - 2\mu_k^x \left[\mu_{q\{1/(\sigma_k^x)^2\}} \sum_{i=1}^n \mu_{q(a_{ik})} \{E_q(\mathbf{x}_{\text{unobs}})(\mathbf{x})\}_i + \frac{\mu_\mu}{\sigma_\mu^2} \right] \right).$$

Therefore,

$q^*(\boldsymbol{\mu})$ is the product of $N(\mu_{q(\mu_k^x)}, \sigma_{q(\mu_k^x)}^2)$ density functions,

$$\text{where } \mu_{q(\mu_k^x)} = \sigma_{q(\mu_k^x)}^2 \left(\mu_{q(1/(\sigma_k^x)^2)} \sum_{i=1}^n \mu_{q(a_{ik})} \{E_q(\mathbf{x}_{\text{unobs}})(\mathbf{x})\}_i + \frac{\mu_\mu}{\sigma_\mu^2} \right), \text{ and}$$

$$\sigma_{q(\mu_k^x)}^2 = 1 / \left(\frac{1}{\sigma_\mu^2} + \mu_{q(1/(\sigma_k^x)^2)} \mu_{q(a_{\bullet k})} \right), \quad 1 \leq k \leq K.$$

Expressions for $A_{\{(\sigma_k^x)^2\}}$, $B_{\{(\sigma_k^x)^2\}}$ and $\mu_{q\{1/(\sigma_k^x)^2\}}$, $1 \leq k \leq K$

To begin, we note that

$$p\{(\boldsymbol{\sigma}^x)^2 | \text{rest}\} \propto p(\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, \boldsymbol{\sigma}^x) p\{(\boldsymbol{\sigma}^x)^2 | \mathbf{a}^x\}$$

$$= \prod_{k=1}^K \prod_{i=1}^n \left[\{2\pi(\sigma_k^x)^2\}^{-1/2} \exp\left\{-\frac{(x_i - \mu_k^x)^2}{2(\sigma_k^x)^2}\right\} \right]^{a_{ik}}.$$

Taking the logarithm of both sides gives

$$\log p\{(\boldsymbol{\sigma}^x)^2 | \text{rest}\} = -\frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n a_{ik} \left[\log\{(\sigma_k^x)^2\} + \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2} \right] - \sum_{k=1}^K \frac{3}{2} \log\{(\sigma_k^x)^2\}$$

$$- \sum_{k=1}^K (1/a_k^x) / (\sigma_k^x)^2 + \text{const}$$

$$= -\sum_{k=1}^K \left\{ \frac{1}{2} (a_{\bullet k} + 1) + 1 \right\} \log\{(\sigma_k^x)^2\}$$

$$- \frac{1}{2} \sum_{k=1}^K \left(a_{\bullet k} \sum_{i=1}^n (x_i - \mu_k^x)^2 + 1/a_k^x \right) / (\sigma_k^x)^2 + \text{const.}$$

Taking expectations:

$$E_q [\log\{(\boldsymbol{\sigma}^x)^2 | \text{rest}\}] = -\sum_{k=1}^K \left\{ \frac{1}{2} (\mu_{q(a_{\bullet k})} + 1) + 1 \right\} \log\{(\sigma_k^x)^2\}$$

$$- \frac{1}{2} \sum_{k=1}^K \left(\mu_{q(a_{\bullet k})} \sum_{i=1}^n E_q \{(x_i - \mu_k^x)^2\} + \mu_{q(1/a_k^x)} \right) / (\sigma_k^x)^2 + \text{const}$$

where

$$E_q \left\{ (x_i - \mu_k^x)^2 \right\} = \left\{ \left(\left\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}) \right\}_i - \mu_{q(\mu_k^x)} \right)^2 + \left\{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \right\}_i + \sigma_{q(\mu_k^x)}^2 \right\}.$$

Therefore,

$q^* \left\{ (\boldsymbol{\sigma}^x)^2 \right\}$ is the product of Inverse-Gamma $\left(A_{q\{(\sigma_k^x)\}}, B_{q\{(\sigma_k^x)\}} \right)$ density functions,

$$\begin{aligned} \text{where } A_{q\{(\sigma_k^x)\}} &= \frac{1}{2} \mu_{q(a_{\bullet k})} + \frac{1}{2}, \\ B_{q\{(\sigma_k^x)\}} &= \mu_{q(1/a_k^x)} + \frac{1}{2} \sum_{i=1}^n \mu_{q(a_{ik})} E_q \left\{ (x_i - \mu_k^x)^2 \right\}, 1 \leq k \leq K. \end{aligned}$$

In addition, using Result 1.4.3,

$$\mu_{q\{1/(\sigma_k^x)^2\}} = A_{q\{(\sigma_k^x)\}} / B_{q\{(\sigma_k^x)\}}.$$

Expressions for $B_{(a_k^x)}$ and $\mu_{q(1/a_k^x)}$, $1 \leq k \leq K$

This derivation is similar to the one for $B_{q(\sigma_k^2)}$ and $\mu_{q(1/\sigma_k^2)}$, and so this gives

$q^*(a_k^x)$ is the product of Inverse-Gamma $\left(1, B_{q(a_k^x)} \right)$ density functions,

where

$$\begin{aligned} B_{q(a_k^x)} &= \mu_{q\{1/(\sigma_k^x)^2\}} + (A_k^x)^{-2}, \text{ and} \\ \mu_{q(1/a_k^x)} &= 1/B_{q(a_k^x)}, 1 \leq k \leq K. \end{aligned}$$

Expressions for $\mu_{q(a_{ik})}$ and ν_{ik} , $1 \leq i \leq n, 1 \leq k \leq K$

Firstly, we note that

$$p(\mathbf{a}|\text{rest}) = \prod_{i=1}^n p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK}|\text{rest}),$$

where each $p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK}|\text{rest})$ has a Multinomial $(1; \omega_1, \dots, \omega_K)$ distribution. Next, we work with each i th component of this distribution.

$$\begin{aligned} p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK}|\text{rest}) &\propto p\left\{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \right\} p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK} | \omega_1, \dots, \omega_K) \\ &= \prod_{k=1}^K \left[\left\{ 2\pi (\sigma_k^x)^2 \right\}^{-1/2} \exp \left\{ -\frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2} \right\} \right]^{a_{ik}} \times \prod_{k=1}^K \frac{\omega_k^{a_{ik}}}{a_{ik}!} \end{aligned}$$

Taking the logarithm of both sides gives

$$\log p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK} | \text{rest}) = \sum_{k=1}^K a_{ik} \left(\log(\omega_k) - \frac{1}{2} \left[\log(2\pi) + \log\{(\sigma_k^x)^2\} + \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2} \right] \right) + \text{const.}$$

Taking expectations, we get

$$E_q \{ \log p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK} | \text{rest}) \} = \sum_{k=1}^K a_{ik} \left(E_q \{ \log(\omega_k) \} - \frac{1}{2} \left[\log(2\pi) + E_q \{ \log(\sigma_k^x)^2 \} + E_q \left\{ \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2} \right\} \right] \right) + \text{const.}$$

Using Results 1.4.3 and 1.4.10, this gives

$$\begin{aligned} & E_q \{ \log p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK} | \text{rest}) \} \\ &= \sum_{k=1}^K a_{ik} \left(\psi(\alpha_{q(\omega_k)}) - \psi \left(\sum_{k=1}^K \alpha_{q(\omega_k)} \right) - \frac{1}{2} \left\{ \log(2\pi) + \log \left(B_{q\{(\sigma_k^x)^2\}} \right) \right. \right. \\ &\quad \left. \left. - \psi \left(A_{q\{(\sigma_k^x)^2\}} \right) + \mu_{q\{1/(\sigma_k^x)^2\}} \left(\left\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}^\top \mathbf{x}) \right\}_i + \left\{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \right\}_i \right. \right. \right. \\ &\quad \left. \left. \left. - 2 \left\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}) \right\}_i \mu_{q(\mu_k^x)} + \mu_{q(\mu_k^x)}^2 + \sigma_{q(\mu_k^x)}^2 \right) \right\} \right) + \text{const} \\ &= \sum_{k=1}^K a_{ik} (\nu_{ik} + \text{const}) + \text{const} \end{aligned}$$

where const does not depend on k (e.g. const includes $-\frac{1}{2} \log(2\pi)$ and $-\psi \left(\sum_{k=1}^K \alpha_{q(\omega_k)} \right)$).

Note that the ‘proportional to’ sign \propto allows us to get rid of ‘const’ below.

$$\begin{aligned} q^*(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK}) &\propto \prod_{k=1}^K (e^{\nu_{ik} + \text{const}})^{a_{ik}} \\ &\propto \prod_{k=1}^K \left(\frac{\exp^{\nu_{ik} + \text{const}}}{\sum_{k'=1}^K e^{\nu_{ik'} + \text{const}}} \right)^{a_{ik}}, \end{aligned}$$

since $\mathcal{C} = \frac{1}{\sum_{k'=1}^K e^{\nu_{ik'} + \text{const}}}$ does not depend on k , and $\prod_{k=1}^K \left(\frac{1}{\mathcal{C}} \right)^{a_{ik}} = \frac{1}{\mathcal{C}}$ given that $\sum_{k=1}^K a_{ik} = 1$.

Hence the cancellation of const. Therefore,

$$q^*(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK}) \text{ is the Multinomial}(1; \mu_{q(a_{i1})}, \dots, \mu_{q(a_{iK})}) \text{ distribution,}$$

where

$$\mu_{q(a_{ik})} = \frac{\exp^{\nu_{ik}}}{\sum_{k'=1}^K e^{\nu_{ik'}}}, \quad 1 \leq k \leq K$$

and

$$\begin{aligned} \nu_{ik} &= \psi(\alpha_{q(\omega_k)}) + \frac{1}{2}\psi\left(A_{q\{(\sigma_k^x)^2\}}\right) - \frac{1}{2}\log\left(B_{q\{(\sigma_k^x)^2\}}\right) \\ &\quad - \frac{1}{2}A_{q\{(\sigma_k^x)^2\}}\left(\{E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}^\top \mathbf{x})\}_i + \{\sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x})\}_i\right) \\ &\quad - 2\{E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x})\}_i \mu_{q(\mu_k^x)} + \mu_{q(\mu_k^x)}^2 + \sigma_{q(\mu_k^x)}^2) / B_{q\{(\sigma_k^x)^2\}}. \end{aligned}$$

Expression for $\alpha_{q(\omega_k)}$, $1 \leq k \leq K$

We begin with

$$\begin{aligned} p(\omega_1, \dots, \omega_K | \text{rest}) &\propto p(\omega_1, \dots, \omega_K) \times \prod_{i=1}^n p(\mathbf{a}_{i1}, \dots, \mathbf{a}_{iK} | \omega_1, \dots, \omega_K) \\ &= \prod_{i=1}^n \prod_{k=1}^K \frac{\omega_k^{a_{ik}}}{a_{ik}!} \times \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\Gamma(\alpha_k)} (\omega_k)^{\alpha_k - 1}. \end{aligned}$$

Taking the logarithm of both sides gives

$$\begin{aligned} \log p(\omega_1, \dots, \omega_K | \text{rest}) &= \sum_{i=1}^n \sum_{k=1}^K (a_{ik} \log(\omega_k) - \log(a_{ik}!)) + \log\left(\Gamma\left(\sum_{k=1}^K \alpha_k\right)\right) \\ &\quad - \sum_{k=1}^K \log(\Gamma(\alpha_k)) + \sum_{k=1}^K (\alpha_k - 1) \log(\omega_k) + \text{const} \\ &= \sum_{i=1}^n \sum_{k=1}^K a_{ik} \log(\omega_k) + \sum_{k=1}^K (\alpha_k - 1) \log(\omega_k) + \text{const}. \end{aligned}$$

Taking expectations, we get

$$E_q \{\log p(\omega_1, \dots, \omega_K | \text{rest})\} = \sum_{i=1}^n \sum_{k=1}^K \mu_{q(a_{ik})} \log(\omega_k) + \sum_{k=1}^K (\alpha_k - 1) \log(\omega_k).$$

Therefore,

$q^*(\omega_1, \dots, \omega_K)$ is the Dirichlet $(\alpha_{q(\omega_1)}, \dots, \alpha_{q(\omega_K)})$ density function,

where

$$\alpha_{q(\omega_k)} = \mu_{q(a_{\cdot k})} + \alpha_k, \quad 1 \leq k \leq K.$$

6.C Derivation of the marginal log-likelihood lower bound

The expression for the lower bound on the marginal log-likelihood given in (6.7) is

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) = & -\frac{1}{2}(K+5)\log(\pi) + \frac{1}{2}(p+q+K+n_{\text{unobs}}) - \frac{p}{2}\log(\sigma_\beta^2) + \frac{1}{2}\log|\boldsymbol{\Sigma}_{q(\beta)}| - \log(A_\varepsilon) \\
 & -\frac{1}{2}\left(2n+n_{\text{unobs}} + \sum_{k=1}^K \sum_{i=1}^n \mu_{q(a_{ik})}\right)\log(2\pi) - \frac{1}{2\sigma_\beta^2}\left\{\|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)})\right\} + \log\Gamma\left\{\frac{1}{2}(n+1)\right\} \\
 & -\frac{1}{2}(n+1)\log(B_{q(\sigma_\varepsilon^2)}) - \log(B_{q(a_\varepsilon)}) + \mu_{q(1/\sigma_\varepsilon^2)}\mu_{q(1/a_\varepsilon)} - \frac{q}{2}\log(\sigma_\alpha^2) + \frac{1}{2}\log|\boldsymbol{\Sigma}_{q(\alpha)}| \\
 & -\log(A_o) - \frac{1}{2\sigma_\alpha^2}\left\{\|\boldsymbol{\mu}_{q(\alpha)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\alpha)})\right\} - \sum_{k=1}^K \log(A_k^x) + \log\Gamma\left\{\frac{1}{2}(n+1)\right\} \\
 & -\frac{K}{2}\log(\sigma_\mu^2) - \frac{1}{2}(n+1)\log(B_{q(\sigma_o^2)}) - \log(B_{q(a_o)}) + \mu_{q(1/\sigma_o^2)}\mu_{q(1/a_o)} \\
 & + \frac{n_{\text{unobs}}}{2}\log(\sigma_{q(\mathbf{x}_{\text{unobs}})}^2) - \frac{1}{2\sigma_\mu^2}\sum_{k=1}^K\left\{\left(\mu_{q(\mu_k^x)} - \mu_\mu\right)^2 + \sigma_{q(\mu_k^x)}^2\right\} + \frac{1}{2}\sum_{k=1}^K \log(\sigma_{q(\mu_k^x)}^2) \\
 & + \sum_{k=1}^K \log\Gamma\left\{\frac{1}{2}(\mu_{q(a_{\bullet k})} + 1)\right\} - \frac{1}{2}\sum_{k=1}^K (\mu_{q(a_{\bullet k})} + 1)\log\left(B_{q\{(\sigma_k^x)^2\}}\right) - \sum_{k=1}^K \log(B_{q(a_k^x)}) \\
 & + \sum_{k=1}^K \mu_{q\{1/(\sigma_k^x)^2\}}\mu_{q(1/a_k^x)} - \sum_{i=1}^n \sum_{k=1}^K \log(\mu_{q(a_{ik})})\mu_{q(a_{ik})} - \log\Gamma\left(\sum_{k=1}^K \alpha_{q(\omega_k)}\right) \\
 & + \log\Gamma\left(\sum_{k=1}^K \alpha_k\right) - \sum_{k=1}^K \log\Gamma(\alpha_k) + \sum_{k=1}^K \log\Gamma(\alpha_{q(\omega_k)}).
 \end{aligned}$$

Derivation: The lower bound on the marginal log-likelihood is achieved through the following expression:

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q) = & E_q\left[\log p\left\{\mathbf{y}, \boldsymbol{\beta}, \sigma_\varepsilon^2, a_\varepsilon, \boldsymbol{o}, \boldsymbol{\alpha}, \sigma_o^2, a_o, \mathbf{x}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2, \mathbf{a}^x, \mathbf{a}, \boldsymbol{\omega}\right\}\right. \\
 & \left. - \log q^*\left\{\boldsymbol{\beta}, \sigma_\varepsilon^2, a_\varepsilon, \boldsymbol{\alpha}, \sigma_o^2, a_o, \mathbf{x}_{\text{unobs}}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2, \mathbf{a}^x, \mathbf{a}, \boldsymbol{\omega}\right\}\right] \\
 = & E_q\left\{\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)\right\} + E_q\left\{\log p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)\right\} \\
 & + E_q\left[\log p\left\{\mathbf{x}_{\text{obs}}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\right\}\right] + E_q\left\{\log p(\boldsymbol{\beta}) - \log q^*(\boldsymbol{\beta})\right\} \\
 & + E_q\left\{\log p(\sigma_\varepsilon^2|a_\varepsilon) - \log q^*(\sigma_\varepsilon^2)\right\} + E_q\left\{\log p(a_\varepsilon) - \log q^*(a_\varepsilon)\right\} \\
 & + E_q\left\{\log p(\boldsymbol{\alpha}) - \log q^*(\boldsymbol{\alpha})\right\} + E_q\left\{\log p(\sigma_o^2|a_o) - \log q^*(\sigma_o^2)\right\} \\
 & + E_q\left\{\log p(a_o) - \log q^*(a_o)\right\} + E_q\left\{\log p(\boldsymbol{\mu}^x) - \log q^*(\boldsymbol{\mu}^x)\right\} \\
 & + E_q\left[\log p\left\{\mathbf{x}_{\text{unobs}}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\right\} - \log q^*(\mathbf{x}_{\text{unobs}})\right] \\
 & + E_q\left[\log p\left\{(\boldsymbol{\sigma}^x)^2|\mathbf{a}^x\right\} - \log q^*\left\{(\boldsymbol{\sigma}^x)^2\right\}\right] + E_q\left\{\log p(\mathbf{a}^x) - \log q^*(\mathbf{a}^x)\right\} \\
 & + E_q\left\{\log \prod_{i=1}^n p(a_{i1}, \dots, a_{iK}|\omega_1, \dots, \omega_K) - \log \prod_{i=1}^n q^*(a_{i1}, \dots, a_{iK})\right\} \\
 & + E_q\left\{\log p(\omega_1, \dots, \omega_K) - \log q^*(\omega_1, \dots, \omega_K)\right\}.
 \end{aligned}$$

6.C. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

First we note that

$$\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

Therefore,

$$E_q \{ \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \} = -\frac{n}{2} \log(2\pi) - \frac{n}{2} E_q \{ \log(\sigma_\varepsilon^2) \} - \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

where

$$\begin{aligned} E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 &= \|\mathbf{y}_{x_{\text{obs}}} - \mathbf{X}_{x_{\text{obs}}} \boldsymbol{\mu}_q(\boldsymbol{\beta})\|^2 + \text{tr} \left(\mathbf{X}_{x_{\text{obs}}}^\top \mathbf{X}_{x_{\text{obs}}} \boldsymbol{\Sigma}_q(\boldsymbol{\beta}) \right) \\ &\quad + \|\mathbf{y}_{x_{\text{unobs}}}\|^2 - 2\mathbf{y}_{x_{\text{unobs}}}^\top E_q(\mathbf{x}_{\text{unobs}}) (\mathbf{X}_{x_{\text{unobs}}}) \boldsymbol{\mu}_q(\boldsymbol{\beta}) \\ &\quad + \text{tr} \left[E_q(\mathbf{x}_{\text{unobs}}) \left(\mathbf{X}_{x_{\text{unobs}}}^\top \mathbf{X}_{x_{\text{unobs}}} \right) \left(\boldsymbol{\Sigma}_q(\boldsymbol{\beta}) + \boldsymbol{\mu}_q(\boldsymbol{\beta}) \boldsymbol{\mu}_q(\boldsymbol{\beta})^\top \right) \right]. \end{aligned}$$

Similarly,

$$E_q \{ \log p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \} = -\frac{n}{2} \log(2\pi) - \frac{n}{2} E_q \{ \log(\sigma_o^2) \} - \frac{1}{2} \mu_{q(1/\sigma_o^2)} E_q \|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2$$

where

$$\begin{aligned} E_q \|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2 &= \|\mathbf{o}_{x_{\text{obs}}} - \tilde{\mathbf{X}}_{x_{\text{obs}}} \boldsymbol{\mu}_q(\boldsymbol{\alpha})\|^2 + \text{tr} \left(\tilde{\mathbf{X}}_{x_{\text{obs}}}^\top \tilde{\mathbf{X}}_{x_{\text{obs}}} \boldsymbol{\Sigma}_q(\boldsymbol{\alpha}) \right) \\ &\quad + \|\mathbf{o}_{x_{\text{unobs}}}\|^2 - 2\mathbf{o}_{x_{\text{unobs}}}^\top E_q(\mathbf{x}_{\text{unobs}}) \left(\tilde{\mathbf{X}}_{x_{\text{unobs}}} \right) \boldsymbol{\mu}_q(\boldsymbol{\alpha}) \\ &\quad + \text{tr} \left[E_q(\mathbf{x}_{\text{unobs}}) \left(\tilde{\mathbf{X}}_{x_{\text{unobs}}}^\top \tilde{\mathbf{X}}_{x_{\text{unobs}}} \right) \left(\boldsymbol{\Sigma}_q(\boldsymbol{\alpha}) + \boldsymbol{\mu}_q(\boldsymbol{\alpha}) \boldsymbol{\mu}_q(\boldsymbol{\alpha})^\top \right) \right]. \end{aligned}$$

Next,

$$\begin{aligned} \log p \{ \mathbf{x}_{\text{obs}} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{obs}}} a_{ik} \left[\log \{ 2\pi (\sigma_k^x)^2 \} + \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2} \right] \\ &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{obs}}} a_{ik} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K a_{\bullet k} \log \{ (\sigma_k^x)^2 \} \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{obs}}} a_{ik} \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2}. \end{aligned}$$

Taking expectations of both sides gives

$$\begin{aligned} E_q \left[\log p \{ \mathbf{x}_{\text{obs}} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \right] &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{obs}}} \mu_{q(a_{ik})} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K \mu_{q(a_{\bullet k})} E_q \left[\log \{ (\sigma_k^x)^2 \} \right] \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^{n_{\text{obs}}} \mu_{q(a_{ik})} \mu_{q\{1/(\sigma_k^x)^2\}} E_q \{ (x_i - \mu_k^x) \} \end{aligned}$$

where

$$E_q \{ (x_i - \mu_k^x)^2 \} = \left\{ \left(x_i - \mu_{q(\mu_k^x)} \right)^2 + \sigma_{q(\mu_k^x)}^2 \right\}.$$

6.C. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

Moving on to the next expectation update, we have

$$\log p(\boldsymbol{\beta}) = -\frac{3}{2} \log(2\pi) - \frac{3}{2} \log(\sigma_\beta^2) - \frac{1}{2\sigma_\beta^2} \|\boldsymbol{\beta}\|^2$$

Taking the expectation of both sides gives:

$$\begin{aligned} E_q \{\log p(\boldsymbol{\beta})\} &= -\frac{p}{2} \log(2\pi) - \frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{2\sigma_\beta^2} E_q \{\|\boldsymbol{\beta}\|^2\} \\ &= -\frac{p}{2} \log(2\pi) - \frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right\}. \end{aligned}$$

The entropy of $q^*(\boldsymbol{\beta})$ is

$$\begin{aligned} -E_q \{\log q^*(\boldsymbol{\beta})\} &= \frac{1}{2} \log \left\{ (2\pi e)^p |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}| \right\} \\ &= \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}| + \frac{p}{2} \log(2\pi) + \frac{p}{2}. \end{aligned}$$

Therefore

$$E_q \{\log p(\boldsymbol{\beta}) - \log q^*(\boldsymbol{\beta})\} = -\frac{p}{2} \log(\sigma_\beta^2) - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right\} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}| + \frac{p}{2}.$$

Next,

$$\begin{aligned} \log p(\sigma_\varepsilon^2 | a_\varepsilon) &= -\frac{1}{2} \log(a_\varepsilon) - \log \Gamma\left(\frac{1}{2}\right) - \frac{3}{2} \log(\sigma_\varepsilon^2) - (1/a_\varepsilon) / \sigma_\varepsilon^2 \\ &= -\frac{1}{2} \log(a_\varepsilon) - \frac{1}{2} \log(\pi) - \frac{3}{2} \log(\sigma_\varepsilon^2) - (1/a_\varepsilon) / \sigma_\varepsilon^2. \end{aligned}$$

Taking expectations, we get

$$E_q \{\log p(\sigma_\varepsilon^2 | a_\varepsilon)\} = -\frac{1}{2} E_q \{\log(a_\varepsilon)\} - \frac{1}{2} \log(\pi) - \frac{3}{2} E_q \{\log(\sigma_\varepsilon^2)\} - \mu_{q(1/a_\varepsilon)} \mu_{q(1/\sigma_\varepsilon^2)}.$$

Also,

$$\begin{aligned} E_q \{\log q^*(\sigma_\varepsilon^2)\} &= \frac{1}{2} (n+1) \log B_{q(\sigma_\varepsilon^2)} - \log \Gamma\left\{\frac{1}{2}(n+1)\right\} - \left\{\frac{1}{2}(n+1) + 1\right\} \log(\sigma_\varepsilon^2) \\ &\quad - B_{q(\sigma_\varepsilon^2)} / \sigma_\varepsilon^2. \end{aligned}$$

Therefore,

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

Next we have

$$\log p(a_\varepsilon) = -\log(A_\varepsilon) - \frac{3}{2} \log(a_\varepsilon) - \frac{1}{2} \log(\pi) - (1/A_\varepsilon^2) / a_\varepsilon.$$

6.C. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

Taking the expectation of both sides gives

$$E_q \{ \log p(a_\varepsilon) \} = -\log(A_\varepsilon) - \frac{3}{2} E_q \{ \log(a_\varepsilon) \} - \frac{1}{2} \log(\pi) - (1/A_\varepsilon^2) \mu_{q(1/a_\varepsilon)}.$$

Also,

$$E_q \{ \log q^*(a_\varepsilon) \} = \log(B_{q(a_\varepsilon)}) - 2E_q \{ \log(a_\varepsilon) \} - B_{q(a_\varepsilon)} \mu_{q(1/a_\varepsilon)}.$$

Therefore,

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

Similarly to $E_q \{ \log p(\boldsymbol{\beta}) - \log q^*(\boldsymbol{\beta}) \}$,

$$E_q \{ \log p(\boldsymbol{\alpha}) - \log q^*(\boldsymbol{\alpha}) \} = -\frac{q}{2} \log(\sigma_\alpha^2) - \frac{1}{2\sigma_\alpha^2} \{ \|\boldsymbol{\mu}_{q(\boldsymbol{\alpha})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\alpha})}) \} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\alpha})}| + \frac{q}{2}.$$

Similarly to $E_q \{ \log p(\sigma_\varepsilon^2|a_\varepsilon) - \log q^*(\sigma_\varepsilon^2) \}$,

$$\begin{aligned} E_q \{ \log p(\sigma_o^2|a_o) - \log q^*(\sigma_o^2) \} &= -\frac{1}{2} E_q \{ \log(a_o) \} + \log \Gamma \left\{ \frac{1}{2}(n+1) \right\} - \frac{1}{2} \log(\pi) \\ &\quad + (B_{q(\sigma_o^2)} - \mu_{q(1/a_o)}) \mu_{q(1/\sigma_o^2)} + \frac{n}{2} E_q \{ \log(\sigma_o^2) \} \\ &\quad - \frac{1}{2}(n+1) \log(B_{q(\sigma_o^2)}). \end{aligned}$$

Similarly to $E_q \{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \}$,

$$\begin{aligned} E_q \{ \log p(a_o) - \log q^*(a_o) \} &= -\log(A_o) - \frac{1}{2} \log(\pi) - \log(B_{q(a_o)}) \\ &\quad + (B_{q(a_o)} - 1/A_o^2) \mu_{q(1/a_o)} + \frac{1}{2} E_q \{ \log(a_o) \}. \end{aligned}$$

Next,

$$\begin{aligned} \log p \{ \mathbf{x}_{\text{unobs}} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n a_{ik} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K a_{\bullet k} \log \{ (\sigma_k^x)^2 \} \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n a_{ik} \frac{(x_i - \mu_k^x)^2}{(\sigma_k^x)^2}. \end{aligned}$$

Taking expectations of both sides gives

$$\begin{aligned} E_q \left[\log p \{ \mathbf{x}_{\text{unobs}} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \right] &= -\frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n \mu_{q(a_{ik})} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K \mu_{q(a_{\bullet k})} E_q \left[\log \{ (\sigma_k^x)^2 \} \right] \\ &\quad - \frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n \mu_{q(a_{ik})} \mu_{q\{1/(\sigma_k^x)^2\}} E_q \{ (x_i - \mu_k^x) \} \end{aligned}$$

where

$$E_q \{ (x_i - \mu_k^x)^2 \} = \left\{ \left(E_q(\mathbf{x}_{\text{unobs}}) \right)_i - \mu_{q(\mu_k^x)} \right\}^2 + \left\{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \right\}_i + \sigma_{q(\mu_k^x)}^2 \}.$$

6.C. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

The entropy of $q^*(\mathbf{x}_{\text{unobs}})$ is

$$-E_q \{\log q^*(\mathbf{x}_{\text{unobs}})\} = \frac{1}{2} \log |\sigma_{q(\mathbf{x}_{\text{unobs}})}^2 \mathbf{I}_{n_{\text{unobs}}}| + \frac{n_{\text{unobs}}}{2} \log(2\pi) + \frac{n_{\text{unobs}}}{2}.$$

Therefore,

$$\begin{aligned} E_q [\log p \{\mathbf{x}_{\text{unobs}} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} - \log q^*(\mathbf{x}_{\text{unobs}})] = \\ -\frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n \mu_{q(a_{ik})} \log(2\pi) - \frac{1}{2} \sum_{k=1}^K \mu_{q(a_{\bullet k})} E_q [\log \{(\sigma_k^x)^2\}] \\ -\frac{1}{2} \sum_{k=1}^K \sum_{i=n_{\text{obs}}+1}^n \mu_{q(a_{ik})} \mu_{q\{1/(\sigma_k^x)^2\}} E_q \{(x_i - \mu_k^x)\} + \frac{n_{\text{unobs}}}{2} \log(\sigma_{q(\mathbf{x}_{\text{unobs}})}^2) \\ + \frac{n_{\text{unobs}}}{2} \log(2\pi) + \frac{n_{\text{unobs}}}{2}. \end{aligned}$$

Moving on, we next have

$$\log p(\boldsymbol{\mu}^x) = -\frac{K}{2} \log(2\pi) - \frac{K}{2} \log(\sigma_\mu^2) - \frac{1}{2\sigma_\mu^2} \sum_{k=1}^K \|\mu_k^x - \mu_\mu\|^2.$$

Taking expectations, we get

$$\begin{aligned} E_q \{\log p(\boldsymbol{\mu}^x)\} &= -\frac{K}{2} \log(2\pi) - \frac{K}{2} \log(\sigma_\mu^2) - \frac{1}{2\sigma_\mu^2} \sum_{k=1}^K E_q \{\|\mu_k^x - \mu_\mu\|^2\} \\ &= -\frac{K}{2} \log(2\pi) - \frac{K}{2} \log(\sigma_\mu^2) - \frac{1}{2\sigma_\mu^2} \sum_{k=1}^K \left\{ \|\mu_{q(\mu_k^x)} - \mu_\mu\|^2 + \sigma_{q(\mu_k^x)} \right\}. \end{aligned}$$

The entropy of $q^*(\boldsymbol{\mu}^x)$ is

$$-E_q \{\log q^*(\boldsymbol{\mu}^x)\} = \frac{K}{2} + \frac{K}{2} \log(2\pi) + \frac{1}{2} \sum_{k=1}^K \log(\sigma_{q(\mu_k^x)}).$$

Therefore,

$$\begin{aligned} E_q \{\log p(\boldsymbol{\mu}^x) - \log q^*(\boldsymbol{\mu}^x)\} &= -\frac{K}{2} \log(\sigma_\mu^2) - \frac{1}{2\sigma_\mu^2} \sum_{k=1}^K \left\{ \|\mu_{q(\mu_k^x)} - \mu_\mu\|^2 + \sigma_{q(\mu_k^x)} \right\} + \frac{K}{2} \\ &\quad + \frac{1}{2} \sum_{k=1}^K \log(\sigma_{q(\mu_k^x)}). \end{aligned}$$

Next,

$$\log p \{(\boldsymbol{\sigma}^x)^2 | \mathbf{a}^x\} = -\frac{1}{2} \sum_{k=1}^K \log(a_k^x) - \frac{K}{2} \log(\pi) - \frac{3}{2} \sum_{k=1}^K \log \{(\sigma_k^x)^2\} - \sum_{k=1}^K (1/a_k^x) / (\sigma_k^x)^2$$

Taking expectations, we get

$$\begin{aligned} E_q [\log p \{(\boldsymbol{\sigma}^x)^2 | \mathbf{a}^x\}] &= -\frac{1}{2} \sum_{k=1}^K E_q \{\log(a_k^x)\} - \frac{K}{2} \log(\pi) - \frac{3}{2} \sum_{k=1}^K E_q [\log \{(\sigma_k^x)^2\}] \\ &\quad - \sum_{k=1}^K \mu_{q(1/a_k^x)} \mu_{q\{1/(\sigma_k^x)^2\}}. \end{aligned}$$

Also,

$$\begin{aligned} E_q [\log q^* \{(\boldsymbol{\sigma}^x)^2\}] &= \frac{1}{2} \sum_{k=1}^K (\mu_{q(a_{\bullet k})} + 1) \log \left(B_{q\{(\sigma_k^x)^2\}} \right) - \sum_{k=1}^K \log \Gamma \left\{ \frac{1}{2} (\mu_{q(a_{\bullet k})} + 1) \right\} \\ &\quad - \sum_{k=1}^K \left\{ \frac{1}{2} (\mu_{q(a_{\bullet k})} + 1) + 1 \right\} E_q [\log \{(\sigma_k^x)^2\}] \\ &\quad - B_{q\{(\sigma_k^x)^2\}} \mu_{q\{1/(\sigma_k^x)^2\}}. \end{aligned}$$

Therefore,

$$\begin{aligned} E_q [\log p \{(\boldsymbol{\sigma}^x)^2 | \mathbf{a}^x\} - \log q^* \{(\boldsymbol{\sigma}^x)^2\}] &= -\frac{1}{2} \sum_{k=1}^K E_q \{ \log(a_k^x) \} - \frac{K}{2} \log(\pi) + \sum_{k=1}^K \log \Gamma \left\{ \frac{1}{2} (\mu_{q(a_{\bullet k})} + 1) \right\} \\ &\quad + \frac{1}{2} \sum_{k=1}^K \mu_{q(a_{\bullet k})} E_q [\log \{(\sigma_k^x)^2\}] - \frac{1}{2} \sum_{k=1}^K (\mu_{q(a_{\bullet k})} + 1) \log \left(B_{q\{(\sigma_k^x)^2\}} \right) \\ &\quad + \sum_{k=1}^K \mu_{q\{1/(\sigma_k^x)^2\}} \left\{ B_{q\{(\sigma_k^x)^2\}} - \mu_{q(1/a_k^x)} \right\}. \end{aligned}$$

Similarly to $E_q \{ \log p(a_\varepsilon) - \log q^*(a_\varepsilon) \}$,

$$\begin{aligned} E_q \{ \log p(\mathbf{a}^x) - \log q^*(\mathbf{a}^x) \} &= -\sum_{k=1}^K \log(A_k^x) - \frac{1}{2} \log(\pi) - \sum_{k=1}^K \log \left(B_{q(a_k^x)} \right) \\ &\quad + \sum_{k=1}^K \left(B_{q(a_k^x)} - 1/(A_k^x)^2 \right) \mu_{q(1/a_k^x)} + \frac{1}{2} \sum_{k=1}^K E_q \{ \log(a_k^x) \}. \end{aligned}$$

Moving on, we have

$$\log \prod_{i=1}^n p(a_{i1}, \dots, a_{iK} | \omega_1, \dots, \omega_K) = \sum_{i=1}^n \sum_{k=1}^K \{ a_{ik} \log(\omega_k) - \log(a_{ik}) \}.$$

Taking expectations of both sides, we get

$$E_q \left\{ \log \prod_{i=1}^n p(a_{i1}, \dots, a_{iK} | \omega_1, \dots, \omega_K) \right\} = \sum_{i=1}^n \sum_{k=1}^K \left[\mu_{q(a_{ik})} E_q \{ \log(\omega_k) \} - E_q \{ \log(a_{ik}) \} \right].$$

Also,

$$E_q \left\{ \log \prod_{i=1}^n q^*(a_{i1}, \dots, a_{iK}) \right\} = \sum_{i=1}^n \sum_{k=1}^K \left[\mu_{q(a_{ik})} \log(\mu_{q(a_{ik})}) - E_q \{ \log(a_{ik}) \} \right].$$

Therefore,

$$\begin{aligned} E_q \left\{ \log \prod_{i=1}^n p(a_{i1}, \dots, a_{iK} | \omega_1, \dots, \omega_K) - \log \prod_{i=1}^n q^*(a_{i1}, \dots, a_{iK}) \right\} &= \sum_{i=1}^n \sum_{k=1}^K \mu_{q(a_{ik})} \left[E_q \{ \log(\omega_k) \} - \log(\mu_{q(a_{ik})}) \right]. \end{aligned}$$

6.C. DERIVATION OF THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

Next,

$$\log p(\omega_1, \dots, \omega_K) = \log \Gamma \left(\sum_{k=1}^K \alpha_k \right) - \sum_{k=1}^K \log \Gamma(\alpha_k) + \sum_{k=1}^K (\alpha_k - 1) \log(\omega_k).$$

Taking the expectation of both sides gives

$$E_q \{ \log p(\omega_1, \dots, \omega_K) \} = \log \Gamma \left(\sum_{k=1}^K \alpha_k \right) - \sum_{k=1}^K \log \Gamma(\alpha_k) + \sum_{k=1}^K (\alpha_k - 1) E_q \{ \log(\omega_k) \}.$$

Also,

$$\begin{aligned} E_q \{ \log q^*(\omega_1, \dots, \omega_K) \} &= \log \Gamma \left(\sum_{k=1}^K \alpha_{q(\omega_k)} \right) - \sum_{k=1}^K \log \Gamma(\alpha_{q(\omega_k)}) \\ &\quad + \sum_{k=1}^K (\alpha_{q(\omega_k)} - 1) E_q \{ \log(\omega_k) \}. \end{aligned}$$

Therefore,

$$\begin{aligned} E_q \{ \log p(\omega_1, \dots, \omega_K) - \log q^*(\omega_1, \dots, \omega_K) \} &= \log \Gamma \left(\sum_{k=1}^K \alpha_k \right) - \log \Gamma \left(\sum_{k=1}^K \alpha_{q(\omega_k)} \right) \\ &\quad - \sum_{k=1}^K \log \Gamma(\alpha_k) + \sum_{k=1}^K \log \Gamma(\alpha_{q(\omega_k)}) \\ &\quad + (\alpha_k - \alpha_{q(\omega_k)}) E_q \{ \log(\omega_k) \}. \end{aligned}$$

Adding all expressions together, we get the lower bound expression given in (6.7). We also note that the following expressions equate to zero.

$$\sum_{k=1}^K (\alpha_k - \alpha_{q(\omega_k)}) E_q \{ \log(\omega_k) \} + \sum_{i=1}^n \sum_{k=1}^K \mu_{q(a_{ik})} E_q \{ \log(\omega_k) \},$$

$$\{ B_{q(\sigma_\varepsilon^2)} - \mu_{q(1/a_\varepsilon)} \} \mu_{q(1/\sigma_\varepsilon^2)} - \frac{1}{2} \mu_{q(1/\sigma_\varepsilon^2)} E_q \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2,$$

$$\{ B_{q(\sigma_o^2)} - \mu_{q(1/a_o)} \} \mu_{q(1/\sigma_o^2)} - \frac{1}{2} \mu_{q(1/\sigma_o^2)} E_q \|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2,$$

and

$$\begin{aligned} &\sum_{k=1}^K \mu_{q\{1/(\sigma_k^x)^2\}} \left\{ B_{q\{(\sigma_k^x)^2\}} - \mu_{q(1/a_k^x)} \right\} \\ &- \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K \mu_{q(a_{ik})} \mu_{q\{1/(\sigma_k^x)^2\}} \left\{ \left(\{ E_{q(\mathbf{x}_{\text{unobs}})}(\mathbf{x}) \}_i - \mu_{q(\mu_k^x)} \right)^2 + \{ \sigma_{q(\mathbf{x}_{\text{unobs}})}^2(\mathbf{x}) \}_i + \sigma_{q(\mu_k^x)}^2 \right\}. \end{aligned}$$

Chapter 7

Alternative approach based on variational message passing

7.1 Introduction

Throughout this thesis, we have dealt with the situation where approximate inference of a possibly large complex graphical model is required. Until now we have shown the details and benefits of using MFVB for fitting these models. An alternative approach, which produces the same approximation but with a different algebraic system, is *variational message passing (VMP)*. The MFVB inference used by Infer.NET is based on the VMP approach. The advantage of VMP compared with MFVB is that the messages are obtained by using only a handful of algebraic rules. This allows more straightforward extension to arbitrarily large models. In this chapter we only derive the VMP alternative algorithms for two models from earlier chapters. Hence, the extension to arbitrarily large models is not covered here.

The essence of VMP is to pass *messages* on a *factor graph*. These two notions are introduced in Section 7.2. Often, these messages are *exponential family* density functions of model parameters and are generally summarised by their natural parameter vector. In addition, the use of natural parameterisations mean that these rules involve simple algebraic manipulations. As will become apparent in section 7.7.1, VMP makes use of the fact that particular forms of messages occur repeatedly, so algebra corresponding to these forms only need to be carried out once.

Even though VMP can be shown to be more efficient than MFVB in terms of per-

forming the necessary algebra, it can also be notationally demanding to outline the rules involved for general graphical models. This is mainly due to the notational and graphical forms contained in Winn & Bishop (2005), Minka (2005) and Minka & Winn (2009) to differ from each other in terms of the underlying rules associated with VMP. We have found it best to work with the description of VMP as given in Minka (2005) and thus have adopted the same notation here. The primary aim of this chapter is to provide a notationally user friendly guide to the general VMP approach, with examples considering formulation and comparison of the algorithms in Chapters 2 and 6 of this thesis.

Section 7.2 gives description of the factor graph representation of Bayesian hierarchical models. Sections 7.3, 7.4 and 7.5 introduce the natural parameter forms, primitives and function definitions corresponding to the density functions that are used within this chapter. The general VMP algorithm is given in Section 7.6 and the reader is provided a walk through using two examples in Section 7.7. Finally, the derivations of the algorithms presented are given in the Appendices that follow.

7.2 Factor graph representation

Factor graphs play a major role in VMP as they are used to illustrate the connection between random variables in a specified model. The aim is to compute messages which are then passed between factors and corresponding nodes on a factor graph. Probabilistic DAGs have direct factor graph representations, e.g., given a DAG that illustrates the conditional independence structure of, say, the random variables x_1, \dots, x_k , these representations arise from relationships such as

$$p(x_1, \dots, x_k) = \prod_{i=1}^k p(x_i | \text{parents of } x_i). \quad (7.1)$$

Another way to represent a model by its factor graph is to derive the joint density function of the random variables and the observed data in terms of their corresponding factors. For example, consider the simple linear regression model

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \sigma^2 &\sim \text{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}), & \boldsymbol{\beta} &\sim \text{N}(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}), \\ \sigma^2 | a &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a), & a &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/A^2). \end{aligned} \quad (7.2)$$

The joint density function of \mathbf{y} , $\boldsymbol{\beta}$, σ^2 and a is

$$p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a) = p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2) p(\sigma^2 | a) p(\boldsymbol{\beta}) p(a), \quad (7.3)$$

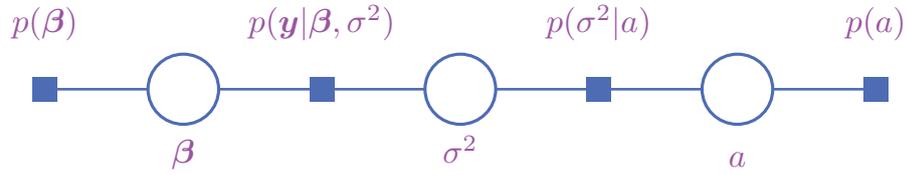


Figure 7.1: Factor graph corresponding to model (7.2).

where each of the density functions $p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$, $p(\sigma^2|a)$, $p(\boldsymbol{\beta})$ and $p(a)$ are referred to as *factors*. The corresponding factor graph is given in Figure 7.1. The circular nodes correspond to each of the random variables in (7.2) and the solid squares correspond to each of the factors in the right-hand side of (7.3), which can also be derived from (7.1). To keep notation simple, we suppress the dependence of $\boldsymbol{\beta}$ and σ^2 on the observed data vector \mathbf{y} in the factor graph. Minka (2005) argues that factor graph representations of models are useful for VMP approximate inference as will be illustrated in the sections to follow.

7.2.1 Additional notation

When using VMP it is convenient to take advantage of the benefits of using a factor graph representation of a statistical model. Factor neighbours notation play an important role in the understanding of VMP in the context of factor graph representations. The following definitions are key to the introduction of Algorithm 10 in Section 7.6.

Definition 7.2.1. Let A and B be sets such that $B \subseteq A$. Then

$$A \setminus B = \text{set containing elements of } A \text{ that are not in } B.$$

We define N_{hid} to be the number of hidden nodes and N_{fac} to be the number of factors in a factor graph. For example, in Figure 7.2, $N_{\text{hid}} = 10$ and $N_{\text{fac}} = 9$ since $\theta_1, \dots, \theta_9$ and ϕ are hidden nodes and f_1, \dots, f_9 are factor nodes.

Definition 7.2.2. Consider a factor graph with factors f_j , $1 \leq j \leq N_{\text{fac}}$ and non-factor nodes θ_i , $1 \leq i \leq N_{\text{hid}}$. Then for each $j = 1, \dots, N_{\text{fac}}$,

$$\text{neighbours}(j) = \{1 \leq i \leq N_{\text{hid}} : \theta_i \text{ is a neighbour of } f_j\}.$$

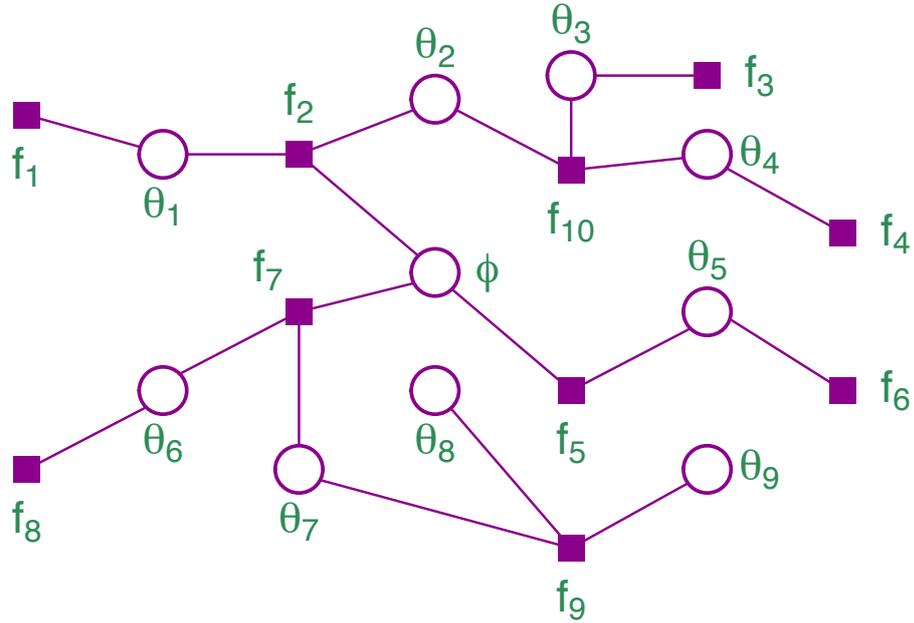


Figure 7.2: *Example of a factor graph.*

It follows that each f_j is a function of its corresponding neighbours on the factor graph.

7.3 Natural parameter forms

Here we provide the forms of the density functions used for the derivations of Algorithms 11 and 12 given later in this chapter.

Univariate normal distribution

The natural parameter form of the univariate normal distribution with mean μ and variance σ^2 is

$$\begin{aligned}
 p(x; \mu, \sigma^2) &= \exp \left\{ \begin{bmatrix} x \\ x^2 \end{bmatrix}^\top \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix} - \left(\frac{\mu}{2\sigma^2} + \frac{\log(\sigma^2)}{2} \right) - \frac{1}{2} \log(2\pi) \right\} \\
 &= \exp \left\{ \mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta} - \left(\frac{\mu}{2\sigma^2} + \frac{\log(\sigma^2)}{2} \right) - \frac{1}{2} \log(2\pi) \right\}
 \end{aligned}$$

where

$$\mathbf{T}(\mathbf{x}) \equiv \begin{bmatrix} x \\ x^2 \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \equiv \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix}$$

are the *natural statistic* and *natural parameter* vectors. Under this representation, we define N_{nat} to be such that

$$x \sim N_{\text{nat}}(\eta_1, \eta_2).$$

Fact 7.3.1. *The ordinary and natural parameters of the univariate normal distribution are mapped between each other through:*

$$\begin{cases} \eta_1 = \mu/\sigma^2 \\ \eta_2 = -1/2\sigma^2 \end{cases} \quad \text{and} \quad \begin{cases} \mu = -\eta_1/2\eta_2 \\ \sigma^2 = -1/2\eta_2. \end{cases}$$

Multivariate normal distribution

The multivariate Normal density function with $n \times 1$ mean vector $\boldsymbol{\mu}$ and $n \times n$ variance-covariance matrix $\boldsymbol{\Sigma}$ is given in Definition 1.4.7. This density function can also be written in its natural parameter form as

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \exp \left\{ \begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} - \frac{1}{2}\boldsymbol{\mu}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} - \frac{1}{2}\log|2\pi\boldsymbol{\Sigma}| \right\} \\ &= \exp \left\{ \mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta} - \frac{1}{2}(\boldsymbol{\mu}^\top\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \log|\boldsymbol{\Sigma}|) - \frac{n}{2}\log(2\pi) \right\} \end{aligned}$$

where

$$\mathbf{T}(\mathbf{x}) \equiv \begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^\top) \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix} \equiv \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}$$

are the *natural statistic* and *natural parameters vectors*. We define $N_{\text{nat,vec}}$ to be such that

$$\mathbf{x} \sim N_{\text{nat,vec}}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2,\text{vec}}).$$

Fact 7.3.2. *The ordinary and natural parameters of the multivariate normal distribution are mapped between each other through the relationship:*

$$\begin{cases} \boldsymbol{\eta}_1 = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ \boldsymbol{\eta}_{2,\text{vec}} = -\frac{1}{2}\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{cases} \quad \text{and} \quad \begin{cases} \boldsymbol{\mu} = -\frac{1}{2}\{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1}\boldsymbol{\eta}_1 \\ \boldsymbol{\Sigma} = -\frac{1}{2}\{\text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}})\}^{-1}. \end{cases}$$

Inverse-Gamma distribution

The Inverse-Gamma density function for a scalar variable x is given in Definition 1.4.9 but can also be written in its natural parameter form as

$$\begin{aligned} p(x; A, B) &= \exp \left\{ \begin{bmatrix} \log x \\ 1/x \end{bmatrix}^\top \begin{bmatrix} -A-1 \\ -B \end{bmatrix} + A \log B - \log \Gamma(A) \right\} \\ &= \exp \{ \mathbf{T}(x)^\top \boldsymbol{\eta} + A \log B - \log \Gamma(A) \} \end{aligned}$$

where

$$\mathbf{T}(x) \equiv \begin{bmatrix} \log x \\ 1/x \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \equiv \begin{bmatrix} -A - 1 \\ -B \end{bmatrix}$$

are the *natural statistic* and *natural parameter* vectors for the Inverse-Gamma distribution. We define IG_{nat} to be such that

$$x \sim \text{IG}_{\text{nat}}(\eta_1, \eta_2).$$

Fact 7.3.3. *The ordinary and natural parameters of the Inverse-Gamma distribution are mapped between each other via:*

$$\begin{cases} \eta_1 = -A - 1 \\ \eta_2 = -B \end{cases} \quad \text{and} \quad \begin{cases} A = -\eta_1 - 1 \\ B = -\eta_2. \end{cases}$$

Inverse-Wishart distribution

The Inverse-Wishart density function with $n \times n$ positive definite scale matrix \mathbf{B} and degrees of freedom a is given in Definition 1.4.3.9 but has another representation using its natural parameter form:

$$\begin{aligned} p(\mathbf{X}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \exp \left\{ \begin{bmatrix} \log |\mathbf{X}| \\ \text{vec}(\mathbf{X}^{-1}) \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2}(a + n + 1) \\ -\frac{1}{2}\text{vec}(\mathbf{B}) \end{bmatrix} - \frac{A}{2} \log |\mathbf{B}| - \log C_{n,a} \right\} \\ &= \exp \left\{ \mathbf{T}(\mathbf{X})^\top \boldsymbol{\eta} - \frac{A}{2} \log |\mathbf{B}| - \log C_{n,a} \right\} \end{aligned}$$

where

$$\mathbf{T}(\mathbf{X}) \equiv \begin{bmatrix} \log |\mathbf{X}| \\ \text{vec}(\mathbf{X}^{-1}) \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \equiv \begin{bmatrix} -\frac{1}{2}(a + n + 1) \\ -\frac{1}{2}\text{vec}(\mathbf{B}) \end{bmatrix}$$

are the *natural statistic* and *natural parameter* vectors for the Inverse-Wishart distribution.

We define IW_{nat} to be such that

$$\mathbf{X} \sim \text{IW}_{\text{nat}}(\eta_1, \boldsymbol{\eta}_2).$$

Fact 7.3.4. *The ordinary and natural parameters of the Inverse-Wishart distribution can be mapped between each other via:*

$$\begin{cases} \eta_1 = -\frac{1}{2}(a + n + 1) \\ \boldsymbol{\eta}_2 = -\frac{1}{2}\text{vec}(\mathbf{B}) \end{cases} \quad \text{and} \quad \begin{cases} A = -2\eta_1 - n - 1 \\ \mathbf{B} = -2\text{vec}^{-1}(\boldsymbol{\eta}_2). \end{cases}$$

Multinomial distribution

The Multinomial density function with 1 independent trial and probability of success $\mathbf{p} = (p_1, \dots, p_k)$ is given in Definition 1.4.13 but has another representation using its natural parameter form:

$$\begin{aligned} p(x_1, \dots, x_k; 1, p_1, \dots, p_k) &= \exp \left\{ \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}^\top \begin{bmatrix} \log(p_1) \\ \vdots \\ \log(p_k) \end{bmatrix} \right\} \\ &= \exp \{ \mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta} \} \end{aligned}$$

where

$$\mathbf{T}(\mathbf{X}) \equiv \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_k \end{bmatrix} \equiv \begin{bmatrix} \log(p_1) \\ \vdots \\ \log(p_k) \end{bmatrix}$$

are the *natural statistic* and *natural parameter* vectors. We define M_{nat} to be such that

$$x_1, \dots, x_k \sim M_{\text{nat}}(1, \eta_1, \dots, \eta_k).$$

Fact 7.3.5. *The ordinary and natural parameters of the Multinomial distribution can be mapped between each other via:*

$$\begin{cases} \eta_1 = \log(p_1) \\ \vdots \\ \eta_k = \log(p_k) \end{cases} \quad \text{and} \quad \begin{cases} p_1 = e^{\eta_1} \\ \vdots \\ p_k = e^{\eta_k}. \end{cases}$$

Dirichlet distribution

The Dirichlet density function with parameters $\alpha_1, \dots, \alpha_k$ is given in Definition 1.4.14. Its natural parameter form is expressed as

$$\begin{aligned} p(x_1, \dots, x_k | \alpha_1, \dots, \alpha_k) &= \exp \left\{ \begin{bmatrix} \log(x_1) \\ \vdots \\ \log(x_k) \end{bmatrix}^\top \begin{bmatrix} \alpha_1 - 1 \\ \vdots \\ \alpha_k - 1 \end{bmatrix} \right. \\ &\quad \left. + \sum_{k=1}^K \log \{ \Gamma(\alpha_k) \} - \log \left\{ \Gamma \left(\sum_{k=1}^K \alpha_k \right) \right\} \right\} \\ &= \exp \left\{ \mathbf{T}(\mathbf{x})^\top \boldsymbol{\eta} + \sum_{k=1}^K \log \{ \Gamma(\alpha_k) \} - \log \left\{ \Gamma \left(\sum_{k=1}^K \alpha_k \right) \right\} \right\} \end{aligned}$$

where

$$\mathbf{T}(\mathbf{x}) \equiv \begin{bmatrix} \log(x_1) \\ \vdots \\ \log(x_k) \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_k \end{bmatrix} \equiv \begin{bmatrix} \alpha_1 - 1 \\ \vdots \\ \alpha_k - 1 \end{bmatrix}$$

are the *natural statistic* and *natural parameter* vectors. We define D_{nat} to be such that

$$x_1, \dots, x_k \sim D_{\text{nat}}(\eta_1, \dots, \eta_k).$$

Fact 7.3.6. *The ordinary and natural parameters of the Dirichlet distribution can be mapped between each other via:*

$$\begin{cases} \eta_1 = \alpha_1 - 1 \\ \vdots \\ \eta_k = \alpha_k - 1 \end{cases} \quad \text{and} \quad \begin{cases} p_1 = \eta_1 + 1 \\ \vdots \\ p_k = \eta_k + 1. \end{cases}$$

7.4 Primitive integrals and results

The following primitives and results follow immediately from the results involving moments of random variables which take on distributions from Section 7.3. These moments aid in the derivation of the algorithms considered in this chapter as they are presented in their natural parameter forms.

We define

$p_{N_{\text{nat}}}\left(x; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}\right)$ to be the $N_{\text{nat}}(\eta_1, \eta_2)$ density function in x ,

$p_{N_{\text{nat,vec}}}\left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,\text{vec}} \end{bmatrix}\right)$ to be the $N_{\text{nat,vec}}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2,\text{vec}})$ density function in \mathbf{x} ,

$p_{IG_{\text{nat}}}\left(x; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}\right)$ to be the $IG_{\text{nat}}(\eta_1, \eta_2)$ density function in x and

$p_{IW_{\text{nat}}}\left(\mathbf{X}; \begin{bmatrix} \eta_1 \\ \boldsymbol{\eta}_2 \end{bmatrix}\right)$ to be the $IW_{\text{nat}}(\eta_1, \boldsymbol{\eta}_2)$ density function in \mathbf{X} .

We now present various primitives and results concerning moments of the density functions given in Section 7.3.

Primitive 7.4.1.

$$\int_{-\infty}^{\infty} x^2 p_{N_{\text{nat}}}\left(x; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}\right) dx = \frac{\eta_1^2 - 2\eta_2}{4\eta_2^2}$$

Primitive 7.4.2.

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x-y)^2 p_{N_{nat}} \left(x, y; \begin{bmatrix} \eta_{x,1} \\ \eta_{x,2} \end{bmatrix}, \begin{bmatrix} \eta_{y,1} \\ \eta_{y,2} \end{bmatrix} \right) dx dy \\ &= \frac{\eta_{x,2}(\eta_{x,1}^2 - 2\eta_{x,2})}{4\eta_{x,2}^3} + \frac{\eta_{y,2}(\eta_{y,1}^2 - 2\eta_{y,2})}{4\eta_{y,2}^3} - \frac{\eta_{x,1}\eta_{y,1}}{2\eta_{x,2}\eta_{y,2}} \end{aligned}$$

where $p_{N_{nat}} \left(x, y; \begin{bmatrix} \eta_{x,1} \\ \eta_{x,2} \end{bmatrix}, \begin{bmatrix} \eta_{y,1} \\ \eta_{y,2} \end{bmatrix} \right)$ is the joint density function of x and y .

Primitive 7.4.3. If \mathbf{x} is a $d \times 1$ vector

$$\begin{aligned} & \int_{\mathbb{R}^d} \mathbf{x}^\top \mathbf{x} p_{N_{nat,vec}} \left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,vec} \end{bmatrix} \right) d\mathbf{x} \\ &= \frac{1}{4} \text{tr} \left(\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} \left[\boldsymbol{\eta}_1 \boldsymbol{\eta}_1^\top \{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} - 2\mathbf{I} \right] \right) \end{aligned}$$

Primitive 7.4.4. For any $n \times 1$ vector \mathbf{a} and $n \times p$ matrix \mathbf{B} ,

$$\begin{aligned} & \int_{\mathbb{R}^p} (\mathbf{a} - \mathbf{B}\mathbf{x})(\mathbf{a} - \mathbf{B}\mathbf{x})^\top p_{N_{nat,vec}} \left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,vec} \end{bmatrix} \right) d\mathbf{x} \\ &= \mathbf{a}\mathbf{a}^\top + \{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} \left[\mathbf{a}\boldsymbol{\eta}_1^\top \mathbf{B}^\top + \frac{1}{4} \mathbf{B}\boldsymbol{\eta}_1 \boldsymbol{\eta}_1^\top \mathbf{B}^\top \{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} - \frac{1}{2} \mathbf{B}\mathbf{B}^\top \right] \end{aligned}$$

Primitive 7.4.5. If $[x_1, \dots, x_d]^\top \sim N_{nat,vec}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2,vec})$, then

$$\begin{aligned} & \int_{\mathbb{R}^d} x_i x_j p_{N_{nat,vec}} \left(\mathbf{x}; \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_{2,vec} \end{bmatrix} \right) d\mathbf{x} \\ &= -\frac{1}{8} \left[\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} \boldsymbol{\eta}_1 \right]_i \left[\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} \boldsymbol{\eta}_1 \right]_j + \left[\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,vec}) \}^{-1} \right]_{ij} \end{aligned}$$

Primitive 7.4.6.

$$\int_0^\infty \frac{1}{x} p_{IG_{nat}} \left(x; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \right) dx = \frac{\eta_1 + 1}{\eta_2}$$

Primitive 7.4.7.

$$\int_0^\infty \log(x) p_{IG_{nat}} \left(x; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \right) dx = \log(-\eta_2) - \psi(-\eta_1 - 1)$$

Primitive 7.4.8.

$$\int_S \mathbf{X}^{-1} p_{IW_{nat}} \left(\mathbf{X}; \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} \right) d\mathbf{X} = \left(\eta_1 + \frac{n+1}{2} \right) \{ \text{vec}^{-1}(\boldsymbol{\eta}_2) \}^{-1}$$

where S is the set of all symmetric positive definite $n \times n$ matrices.

Result 7.4.1. Suppose that $x_1, \dots, x_K \sim M_{nat}(1; \eta_1, \dots, \eta_K)$. Then

$$E(x_k) = e^{\eta_k}, \quad 1 \leq k \leq K.$$

Result 7.4.2. Suppose that $x_1, \dots, x_K \sim D_{nat}(\eta_1, \dots, \eta_K)$. Then

$$E \{ \log(x_k) \} = \psi(\eta_k + 1) - \psi \left\{ \sum_{k=1}^K (\eta_k + 1) \right\}, \quad 1 \leq k \leq K.$$

7.5 Function definitions

The following functions are useful for describing various VMP updates used throughout this chapter.

$$\mathcal{F} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \right) \equiv \frac{1}{4} \text{tr} \left(\{ \text{vec}^{-1}(a_2) \}^{-1} \left[a_1 a_1^\top \{ \text{vec}^{-1}(a_2) \}^{-1} - 2\mathbf{I} \right] \right),$$

where a_1 and a_2 are scalars.

$$\mathcal{G} \left(\begin{bmatrix} a_1 \\ \mathbf{a}_2 \end{bmatrix}; \mathbf{b}, \mathbf{C} \right) \equiv \left(a_1 + \frac{n+1}{2} \right) \{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1},$$

where a_1 is a scalar and \mathbf{a}_2 is an $n^2 \times 1$ vector.

$$\mathcal{H} \left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; \mathbf{b}, \mathbf{C} \right) \equiv \mathbf{b}\mathbf{b}^\top + \{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1} \left[\mathbf{b}\mathbf{a}_1^\top \mathbf{C}^\top + \frac{1}{4} \mathbf{C}\mathbf{a}_1\mathbf{a}_1^\top \mathbf{C}^\top \{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1} - \frac{1}{2} \mathbf{C}\mathbf{C}^\top \right],$$

where \mathbf{a}_1 is a $q \times 1$ vector, \mathbf{a}_2 a $q^2 \times 1$ vector, \mathbf{b} a $p \times 1$ vector and \mathbf{C} a $p \times q$ matrix.

$$\mathcal{I} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \right) \equiv \frac{a_2(a_1^2 - 2a_2)}{4a_2^3} + \frac{b_2(b_1^2 - 2b_2)}{4b_2^3} - \frac{a_1 b_1}{2a_2 b_2},$$

where a_1, a_2, b_1 and b_2 are scalars.

$$\mathcal{J} \left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}; b, c \right) \equiv -\frac{1}{8} \left[\{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1} \mathbf{a}_1 \right]_b \left[\{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1} \mathbf{a}_1 \right]_c + \left[\{ \text{vec}^{-1}(\mathbf{a}_2) \}^{-1} \right]_{bc},$$

where \mathbf{a}_1 is a $q \times 1$ vector, \mathbf{a}_2 is a $q^2 \times 1$ vector and b and c both scalars.

$$\mathcal{K} \left(\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \right) \equiv \frac{a_1^2 - 2a_2}{4a_2^2},$$

where a_1 and a_2 are scalars.

7.6 The variational message passing algorithm

Our goal is to obtain the mean field approximation to the joint posterior density function:

$$p(\theta_1, \dots, \theta_{N_{\text{hid}}} | \mathbf{x}) \approx q(\theta_1) \dots q(\theta_{N_{\text{hid}}}) \quad (7.4)$$

where $\theta_1, \dots, \theta_{N_{\text{hid}}}$ are the parameters or hidden nodes and \mathbf{x} contains the data or evidence nodes in the corresponding statistical model. VMP aims to obtain the optimal q -densities through an iterative scheme concerned with updating messages of the form:

$$m_{f_j \rightarrow \theta_i} \quad \text{and} \quad m_{\theta_i \rightarrow f_j} \quad \text{for each } i \in \text{neighbours}(j), \quad 1 \leq j \leq N_{\text{fac}}.$$

The general iterative VMP scheme, labelled Algorithm 10, can be used to solve the optimal q densities, by updating such messages.

Convergence of Algorithm 10 is assessed by monitoring successive values of the lower bound on the marginal log-likelihood $\log \underline{p}(\mathbf{x}; q)$ and stopping once the successive differences are very small. The lower bound on the marginal log-likelihood is given by

$$\log \underline{p}(\mathbf{x}; q) = E_q \left[\sum_{j=1}^{N_{\text{fac}}} \log f_j(\theta_{\text{neighbours}(j)}) - \sum_{i=1}^{N_{\text{hid}}} \log q(\theta_i) \right]. \quad (7.5)$$

Within Algorithm 10, if $\text{neighbours}(j) \setminus \{i\} = \emptyset$, the empty set, the expression

$$\prod_{i' \in \text{neighbours}(j) \setminus \{i\}} m_{\theta_{i'} \rightarrow f_j}(\theta_{i'}) m_{f_j \rightarrow \theta_{i'}}(\theta_{i'})$$

is taken to be 1, since there would exist no i' . Integration is replaced by summation if $\theta_{\text{neighbours}(j) \setminus \{i\}}$ in Step 2 (a) contains discrete random variables. It is easier to work with messages appearing under the same product sign to be from the same exponential family. This is referred to as conjugate VMP. However, VMP is a general scheme and not necessarily restricted to conjugate situations. When conjugacy exists within the message updates, the functional forms of these messages are exposed and it becomes clear that the parameters of these messages are inter-related, much like what happens with MFVB.

Step 1: Initialise

$m_{f_j \rightarrow \theta_i}(\theta_i)$ for all $i \in \text{neighbours}(j)$, $1 \leq j \leq N_{\text{fac}}$, to be arbitrary density functions.

Step 2: Cycle

(a) For $i \in \text{neighbours}(j)$, $1 \leq j \leq N_{\text{fac}}$,

$$m_{\theta_i \rightarrow f_j}(\theta_i) \leftarrow \prod_{\{j' \neq j : i \in \text{neighbours}(j')\}} m_{f_{j'} \rightarrow \theta_i}(\theta_i),$$

(b) For $i \in \text{neighbours}(j)$, $1 \leq j \leq N_{\text{fac}}$,

$$m_{f_j \rightarrow \theta_i}(\theta_i) \leftarrow \exp \left[\int \left(\frac{1}{Z} \prod_{i' \in \text{neighbours}(j) \setminus \{i\}} m_{\theta_{i'} \rightarrow f_j}(\theta_{i'}) m_{f_j \rightarrow \theta_{i'}}(\theta_{i'}) \right) \times \log f_j(\theta_{\text{neighbours}(j)}) d\theta_{\text{neighbours}(j) \setminus \{i\}} \right],$$

$$\text{where } Z = \int \prod_{i' \in \text{neighbours}(j) \setminus \{i\}} m_{\theta_{i'} \rightarrow f_j}(\theta_{i'}) m_{f_j \rightarrow \theta_{i'}}(\theta_{i'}) d\theta_{\text{neighbours}(j) \setminus \{i\}}.$$

until all messages converge.

Step 3: For $1 \leq i \leq N_{\text{hid}}$,

$$q^*(\theta_i) \leftarrow \prod_{\{j : i \in \text{neighbours}(j)\}} m_{f_j \rightarrow \theta_i}^*(\theta_i) / \int \prod_{\{j : i \in \text{neighbours}(j)\}} m_{f_j \rightarrow \theta_i}^*(\theta_i) d\theta_i.$$

Algorithm 10: *The general variational message passing algorithm.*

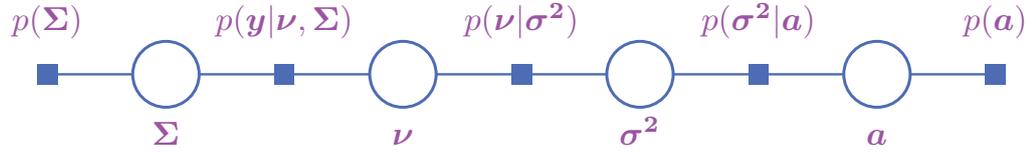


Figure 7.3: Factor graph corresponding to the marginal longitudinal semiparametric regression model in (2.15).

These updates begin to exhibit particular patterns and *remembering* these patterns allows one to streamline the calculations of VMP. This is especially useful for larger complex statistical models. Examples of these patterns are given in Appendix 7.4 for the two models considered in the following sections. In addition, when computing the VMP updates in Algorithm 10, it is appropriate to work with the natural parameter forms of exponential family density functions.

7.7 Examples

We consider two regression examples used throughout this thesis, namely the marginal longitudinal semiparametric regression model (2.15) given in Chapter 2 and the mixture model with measurement error (6.4) given in Chapter 6. We provide detail on the VMP methodology specifically for these two models.

7.7.1 Marginal longitudinal semiparametric regression model

The factor graph corresponding to model (2.15) is given in Figure 7.3, where the circular nodes correspond to each of the random variables, vectors or matrices and the solid squares represent each of the factors in the model. To make notation simple we have combined the mean function coefficients into a single vector, i.e.,

$$\boldsymbol{\nu} \equiv \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}.$$

We aim to obtain a mean field approximation to the joint posterior density function:

$$p(\boldsymbol{\nu}, \mathbf{a}, \sigma^2, \boldsymbol{\Sigma} | \mathbf{y}) \approx q(\boldsymbol{\nu})q(\mathbf{a})q(\sigma^2)q(\boldsymbol{\Sigma}).$$

Following the rules of Algorithm 10 and referring to the factor graph in Figure 7.3 for the stochastic node $\boldsymbol{\Sigma}$, the two factor to node messages to initialise are $\mathbf{m}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$ and $\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$. We can set these to be any density functions, however to make the

calculations simple and impose conjugacy we choose them to be the Inverse-Wishart(n, \mathbf{I}_n) density function. This can be written as:

$$\begin{aligned} \mathbf{m}_{p(\Sigma) \rightarrow \Sigma}(\Sigma) &\leftarrow \exp \left\{ \begin{bmatrix} \log |\Sigma| \\ \text{vec}(\Sigma^{-1}) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\Sigma) \rightarrow \Sigma} \right\} \\ \mathbf{m}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma}(\Sigma) &\leftarrow \exp \left\{ \begin{bmatrix} \log |\Sigma| \\ \text{vec}(\Sigma^{-1}) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma} \right\} \end{aligned}$$

where the natural parameter vectors $\boldsymbol{\eta}_{p(\Sigma) \rightarrow \Sigma}$ and $\boldsymbol{\eta}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma}$ are initialised to be

$$\boldsymbol{\eta}_{p(\Sigma) \rightarrow \Sigma} \leftarrow \begin{bmatrix} -n-1 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_n) \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma} \leftarrow \begin{bmatrix} -n-1 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_n) \end{bmatrix}.$$

The natural parameter forms of the density functions used for this model were given in Section 7.3. Step 2 (a) of Algorithm 10 is concerned with updating the stochastic node to factor messages. If we begin with the message $\mathbf{m}_{\Sigma \rightarrow p(\Sigma)}(\Sigma)$, we notice that the only other neighbour of Σ apart from $p(\Sigma)$ is $p(\mathbf{y}|\nu, \Sigma)$. This means that $\mathbf{m}_{\Sigma \rightarrow p(\Sigma)}(\Sigma)$ has the update

$$\begin{aligned} \mathbf{m}_{\Sigma \rightarrow p(\Sigma)}(\Sigma) &\leftarrow \mathbf{m}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma}(\Sigma) \\ &= \exp \left\{ \begin{bmatrix} \log |\Sigma| \\ \text{vec}(\Sigma^{-1}) \end{bmatrix}^\top \boldsymbol{\eta}_{\Sigma \rightarrow p(\Sigma)} \right\} \end{aligned}$$

where

$$\boldsymbol{\eta}_{\Sigma \rightarrow p(\Sigma)} \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma}.$$

Similarly, $\mathbf{m}_{\Sigma \rightarrow p(\mathbf{y}|\nu, \Sigma)}(\Sigma) \leftarrow \mathbf{m}_{p(\Sigma) \rightarrow \Sigma}(\Sigma)$, which corresponds to the natural parameter update:

$$\boldsymbol{\eta}_{\Sigma \rightarrow p(\mathbf{y}|\nu, \sigma^2)} \leftarrow \boldsymbol{\eta}_{p(\Sigma) \rightarrow \Sigma}.$$

Step 2 (b) of Algorithm 10 is concerned with updating the factor to stochastic node messages and focusing on stochastic node Σ which involves updating $\mathbf{m}_{p(\Sigma) \rightarrow \Sigma}(\Sigma)$ and $\mathbf{m}_{p(\mathbf{y}|\nu, \Sigma) \rightarrow \Sigma}(\Sigma)$. The first factor to stochastic node message to update is $\mathbf{m}_{p(\Sigma) \rightarrow \Sigma}(\Sigma)$. Since the only neighbour of factor $p(\Sigma)$ is the stochastic node Σ this message has the fol-

lowing form

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) &\leftarrow \exp\{\log p(\boldsymbol{\Sigma})\} \\ &\propto \exp\left\{\begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2}(A_\Sigma + n + 1) \\ -\frac{1}{2}\text{vec}(\mathbf{B}_\Sigma) \end{bmatrix}\right\}. \end{aligned}$$

This means that

$$\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} \leftarrow \begin{bmatrix} -\frac{1}{2}(A_\Sigma + n + 1) \\ -\frac{1}{2}\text{vec}(\mathbf{B}_\Sigma) \end{bmatrix}$$

and remains fixed at this value as it doesn't depend on any other messages. The next message has the form:

$$\begin{aligned} \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) &\leftarrow \exp\left\{\int_{\mathbb{R}^v} \mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})}(\boldsymbol{\nu}) \times \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \right. \\ &\quad \left. \times \log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) d(\boldsymbol{\nu})\right\}, \end{aligned}$$

where $v = p + \sum_{\ell=1}^d K_\ell$, the number of columns in $\mathbf{C} = [\mathbf{X} | \mathbf{Z}]$. Further breakdown of $\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$ gives

$$\begin{aligned} &\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) \\ &\propto \exp\left\{\int_{\mathbb{R}^v} \left\{\begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}\right)\right\} \right. \\ &\quad \left. \times \left[\begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^\top \begin{bmatrix} -\frac{m}{2} \\ -\frac{1}{2}\sum_{i=1}^m \text{vec}\left\{(\mathbf{y}_i - \mathbf{C}_i\boldsymbol{\nu}_i)(\mathbf{y}_i - \mathbf{C}_i\boldsymbol{\nu}_i)^\top\right\}\right]\right\} d\boldsymbol{\nu} \Bigg] \\ &\propto \exp\left\{\begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^\top \begin{bmatrix} -\frac{m}{2} \\ -\frac{1}{2}\sum_{i=1}^m \text{vec}\left\{\int_{\mathbb{R}^v} (\mathbf{y}_i - \mathbf{C}_i\boldsymbol{\nu}_i)(\mathbf{y}_i - \mathbf{C}_i\boldsymbol{\nu}_i)^\top \right. \right. \\ \quad \left. \left. \times p_{\text{N}_{\text{nat}, \text{vec}}}(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} \right. \right. \\ \quad \left. \left. + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}) d\boldsymbol{\nu}\right\}\right] \Bigg\} \end{aligned}$$

where $p_{\text{N}_{\text{nat}, \text{vec}}}(\mathbf{x}; [\boldsymbol{\eta}_1^\top \boldsymbol{\eta}_{2, \text{vec}}^\top]^\top)$ is the $\text{N}_{\text{nat}, \text{vec}}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_{2, \text{vec}})$ density function in \mathbf{x} . Further detail on this notation is given previously in Sections 7.3 and 7.4. Using Primitive 7.4.4

we get

$$\begin{aligned} & \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) \\ & \propto \exp \left\{ \begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^\top \begin{bmatrix} -\frac{m}{2} \\ -\frac{1}{2} \sum_{i=1}^m \text{vec} \left\{ \mathcal{H} \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}; \mathbf{y}_i, \mathbf{C}_i} \right) \right\} \end{bmatrix} \right\}, \end{aligned}$$

where the function \mathcal{H} is defined in Section 7.5. This corresponds to the natural parameter update having the form

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} \leftarrow \begin{bmatrix} -\frac{m}{2} \\ -\frac{1}{2} \sum_{i=1}^m \text{vec} \left\{ \mathcal{H} \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}; \mathbf{y}_i, \mathbf{C}_i} \right) \right\} \end{bmatrix}.$$

Step 2 (b) is now finished for the updates corresponding to the stochastic node $\boldsymbol{\Sigma}$. Once convergence is achieved for the messages in Steps 2 (a) and (b) corresponding to all stochastic nodes in factor graph (7.3), we move on to Step 3 of Algorithm 10. For $\boldsymbol{\Sigma}$ this involves the update:

$$\boldsymbol{\eta}_{q(\boldsymbol{\Sigma})}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow (\boldsymbol{\nu})}^* + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow (\boldsymbol{\nu})}^*.$$

This means that

$$q^*(\boldsymbol{\Sigma}) \text{ is the } \text{IW}_{\text{nat}} \left(\boldsymbol{\eta}_{q(\boldsymbol{\Sigma}),1}^*, \boldsymbol{\eta}_{q(\boldsymbol{\Sigma}),2}^* \right) \text{ density function in } \boldsymbol{\Sigma}$$

and can also be transformed back to the original parameters, such as those used in Chapter 2, by writing this as the

$$\text{Inverse-Wishart} \left\{ - \left(2\boldsymbol{\eta}_{q(\boldsymbol{\Sigma}),2}^* + n + 1 \right), -2\text{vec} \left(\boldsymbol{\eta}_{q(\boldsymbol{\Sigma}),1}^* \right) \right\} \text{ density function in } \boldsymbol{\Sigma}.$$

Algorithm 11 gives the updates for all stochastic nodes in Figure 7.3 and the derivation is given in Appendix 7.A. The lower bound on the marginal log-likelihood for this example is derived by using the expression given in (7.5) and gives the same result as in expression (2.17) in Chapter 2. Therefore, once the natural parameters are transformed back to the original parameters, convergence is assessed in the same manner as before.

Initialise:

$$\begin{aligned}
 \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \\
 \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} &\leftarrow \begin{bmatrix} \mathbf{0}_v \\ -\frac{1}{2}\text{vec}(\mathbf{I}_v) \end{bmatrix}; \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} &\leftarrow \begin{bmatrix} \mathbf{0}_v \\ -\frac{1}{2}\text{vec}(\mathbf{I}_v) \end{bmatrix}; & \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} &\leftarrow \begin{bmatrix} -n-1 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_n) \end{bmatrix}; \\
 \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} &\leftarrow \begin{bmatrix} -n-1 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_n) \end{bmatrix}
 \end{aligned}$$

Cycle:

$$\begin{aligned}
 \text{(a)} \quad \boldsymbol{\eta}_{a_\ell \rightarrow p(a_\ell)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \quad ; \quad \boldsymbol{\eta}_{a_\ell \rightarrow p(\sigma_\ell^2|a_\ell)} \leftarrow \boldsymbol{\eta}_{p(a) \rightarrow a} \\
 \boldsymbol{\eta}_{\sigma_\ell^2 \rightarrow p(\sigma_\ell^2|a_\ell)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} \quad ; \quad \boldsymbol{\eta}_{\sigma_\ell^2 \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)} \leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} \\
 \boldsymbol{\eta}_{\boldsymbol{\nu} \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} \quad ; \quad \boldsymbol{\eta}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\sigma}^2)} \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} \\
 \boldsymbol{\eta}_{\boldsymbol{\Sigma} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\sigma}^2)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} \quad ; \quad \boldsymbol{\eta}_{\boldsymbol{\Sigma} \rightarrow p(\boldsymbol{\Sigma})} \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}
 \end{aligned}$$

$$\begin{aligned}
 \text{(b)} \quad \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} &\leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\frac{(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2})_1 + 1}{(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2})_2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} &\leftarrow \begin{bmatrix} -\frac{3}{2} \\ -\frac{(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell})_1 + 1}{(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell})_2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} &\leftarrow \begin{bmatrix} -\frac{1}{2}K_\ell \\ -\frac{1}{2}\mathcal{F}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_\ell} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \mathbf{u}_\ell}) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} &\leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2}\text{vec} \left(\begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}_{1 \leq \ell \leq d} \left\{ \left(\frac{(\boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2})_1 + 1}{(\boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2})_2} \right) \mathbf{I}_{K_\ell} \right\} \end{bmatrix} \right) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} &\leftarrow \begin{bmatrix} \mathbf{C}^\top \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}; n) \right\} \mathbf{y} \\ \frac{1}{2}\text{vec} \left[\mathbf{C}^\top \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}; n) \right\} \mathbf{C} \right] \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} &\leftarrow \begin{bmatrix} -\frac{m}{2} \\ -\frac{1}{2} \sum_{i=1}^m \text{vec} \left\{ \mathcal{H}(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}; \mathbf{y}_i, \mathbf{C}_i) \right\} \end{bmatrix}
 \end{aligned}$$

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

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} \leftarrow \left[\begin{array}{c} -\frac{m}{2} \\ -\frac{1}{2} \sum_{i=1}^m \text{vec} \left\{ \mathcal{H} \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} ; \mathbf{y}_i, \mathbf{C}_i \right) \right\} \end{array} \right]$$

$$\boldsymbol{\eta}_{q(\boldsymbol{\nu})}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow (\boldsymbol{\nu})}^* + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow (\boldsymbol{\nu})}^*$$

$$\boldsymbol{\eta}_{q(\boldsymbol{\Sigma})}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}^* + \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}^*$$

For $1 \leq \ell \leq d$:

$$\boldsymbol{\eta}_{q(a_\ell)}^* \leftarrow \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell}^* + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell}^*$$

$$\boldsymbol{\eta}_{q(\sigma_\ell^2)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}^* + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2}^*$$

Algorithm 11: VMP algorithm for determining the natural parameter vectors $\boldsymbol{\eta}_{q(a_\ell)}^*$, $\boldsymbol{\eta}_{q(\sigma_\ell^2)}^*$, $\boldsymbol{\eta}_{q(\boldsymbol{\nu})}^*$ and $\boldsymbol{\eta}_{q(\boldsymbol{\Sigma})}^*$ of the optimal density functions $q^*(\boldsymbol{\nu})$, $q^*(\boldsymbol{\sigma}^2)$, $q^*(\mathbf{a})$ and $q^*(\boldsymbol{\Sigma})$ for approximate inference in Model (2.15).

7.7.2 Mixture model with measurement error

The factor graph corresponding to model (6.4) is given in Figure 7.4. We impose the same mean field approximations on the joint posterior density functions as given in Chapter 6. Algorithm 12 gives the VMP updates corresponding to this model. Details of the derivations for the updates given in Algorithm 12 are given in Appendix 7.B. In order to impose conjugacy, we initialise the messages corresponding to stochastic nodes \mathbf{a}^x , $(\boldsymbol{\sigma}^x)^2$, σ_ε^2 , a_ε , σ_o^2 and a_o to be the Inverse-Gamma(1,1) density function, corresponding to the natural parameter vector being set to $[-2 \ -1]^\top$. The messages corresponding to the stochastic nodes $\mathbf{x}_{\text{unobs}}$ and $\boldsymbol{\mu}^x$ are initialised to be the $N(0,1)$ density function, which corresponds to the natural parameter vector being set to $[0 \ -\frac{1}{2}]^\top$. The messages relating to the columns in node \mathbf{a} are initialised to be the Multinomial $(a_{i1}, \dots, a_{iK}; 1, \frac{1}{K}, \dots, \frac{1}{K})$ density function, where K is the number of mixture components used. This corresponds to the natural parameters being set to $\log(1/K), \dots, \log(1/K)$. The messages relating to $\boldsymbol{\omega}$ are initialised to be the Dirichlet $(\omega_1, \dots, \omega_K; 2, \dots, 2)$ density function and this corresponds to the natural parameters being set to $1, \dots, 1$. Lastly, the messages corresponding to $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are initialised to be the $N(\mathbf{0}_d, \mathbf{I}_d)$ density function, where d is the dimension of the corresponding stochastic vector node, and this corresponds to the natural parameter vector being set to $[\mathbf{0}_d \ -\frac{1}{2} \text{vec}(\mathbf{I}_d)]^\top$.

Initialise:

$$\begin{aligned}
 \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow (\sigma_k^x)^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} &\leftarrow 1 \\
 \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}} &\leftarrow \log(1/K); & \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow x_{\text{unobs},i}} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; & \boldsymbol{\eta}_{p(a_{ik} | \omega_k) \rightarrow a_{ik}} &\leftarrow \log(1/K) \\
 \boldsymbol{\eta}_{p(\mathbf{a}_k | \omega_k) \rightarrow \omega_k} &\leftarrow 1; & \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} &\leftarrow \begin{bmatrix} \mathbf{0}_2 \\ -\frac{1}{2} \text{vec}(\mathbf{I}_2) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow a_o} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; & \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow \sigma_o^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} &\leftarrow \begin{bmatrix} \mathbf{0}_2 \\ -\frac{1}{2} \text{vec}(\mathbf{I}_2) \end{bmatrix}; & \boldsymbol{\eta}_{p(a_o) \rightarrow a_o} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} &\leftarrow \begin{bmatrix} \mathbf{0}_3 \\ -\frac{1}{2} \text{vec}(\mathbf{I}_3) \end{bmatrix}; & \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} &\leftarrow \begin{bmatrix} -3/2 \\ -1/A_\varepsilon^2 \end{bmatrix}; & \boldsymbol{\eta}_{p(\sigma_\varepsilon^2 | a_\varepsilon) \rightarrow a_\varepsilon} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} &\leftarrow \begin{bmatrix} \mathbf{0}_3 \\ -\frac{1}{2} \text{vec}(\mathbf{I}_3) \end{bmatrix}; & \boldsymbol{\eta}_{p(\sigma_\varepsilon^2 | a_\varepsilon) \rightarrow \sigma_\varepsilon^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}
 \end{aligned}$$

Cycle:

$$\begin{aligned}
 \text{(a)} \quad \boldsymbol{\eta}_{a_k^x \rightarrow p(a_k^x)} &\leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x}; & \boldsymbol{\eta}_{a_k^x \rightarrow p\{(\sigma_k^x)^2 | a_k^x\}} &\leftarrow \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \\
 \boldsymbol{\eta}_{(\sigma_k^x)^2 \rightarrow p\{(\sigma_k^x)^2 | a_k^x\}} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2} \\
 \boldsymbol{\eta}_{(\sigma_k^x)^2 \rightarrow p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow (\sigma_k^x)^2} \\
 \boldsymbol{\eta}_{a_{ik} \rightarrow p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p(a_{ik} | \omega_k) \rightarrow a_{ik}}; & \boldsymbol{\eta}_{a_{ik} \rightarrow p(a_{ik} | \omega_k)} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}} \\
 \boldsymbol{\eta}_{\omega_k \rightarrow p(a_{ik} | \omega_k)} &\leftarrow \boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k}; & \boldsymbol{\eta}_{\omega_k \rightarrow p(\omega_k)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{a}_k | \omega_k) \rightarrow \omega_k} \\
 \boldsymbol{\eta}_{\mu_k^x \rightarrow p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x}; & \boldsymbol{\eta}_{\mu_k^x \rightarrow p(\mu_k^x)} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x} \\
 \boldsymbol{\eta}_{\boldsymbol{\alpha} \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}; & \boldsymbol{\eta}_{\boldsymbol{\alpha} \rightarrow p(\boldsymbol{\alpha})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \\
 \boldsymbol{\eta}_{\sigma_o^2 \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow \sigma_o^2}; & \boldsymbol{\eta}_{\sigma_o^2 \rightarrow p(\sigma_o^2 | a_o)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \\
 \boldsymbol{\eta}_{a_o \rightarrow p(\sigma_o^2 | a_o)} &\leftarrow \boldsymbol{\eta}_{p(a_o) \rightarrow a_o}; & \boldsymbol{\eta}_{a_o \rightarrow p(a_o)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow a_o} \\
 \boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}}
 \end{aligned}$$

7.7. EXAMPLES

$$\boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \leftarrow \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p(o|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}$$

$$\boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p(o|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} \leftarrow \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}$$

$$\boldsymbol{\eta}_{\boldsymbol{\beta} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \leftarrow \boldsymbol{\eta}_p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}; \quad \boldsymbol{\eta}_{\boldsymbol{\beta} \rightarrow p(\boldsymbol{\beta})} \leftarrow \boldsymbol{\eta}_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}$$

$$\boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \leftarrow \boldsymbol{\eta}_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2; \quad \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2$$

$$\boldsymbol{\eta}_{a_\varepsilon \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} \leftarrow \boldsymbol{\eta}_p(a_\varepsilon) \rightarrow a_\varepsilon; \quad \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(a_\varepsilon)} \leftarrow \boldsymbol{\eta}_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon$$

$$(b) \quad \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \leftarrow \begin{bmatrix} -3/2 \\ -1/(A_k^x)^2 \end{bmatrix}; \quad \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \leftarrow \begin{bmatrix} \mu_\mu \sigma_\mu^2 \\ -1/2\sigma_\mu^2 \end{bmatrix}$$

$$\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x \leftarrow \begin{bmatrix} -\frac{1}{2} \\ \frac{(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2)_1}{(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2)_2} + 1 \end{bmatrix}$$

$$\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 \leftarrow \begin{bmatrix} -3/2 \\ \frac{(\boldsymbol{\eta}_p(a_k^x) \rightarrow a_k^x + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x)_1}{(\boldsymbol{\eta}_p(a_k^x) \rightarrow a_k^x + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x)_2} + 1 \end{bmatrix}$$

$$\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}$$

$$\leftarrow \begin{bmatrix} -\frac{1}{2} \exp\left(\boldsymbol{\eta}_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}\right) \\ \times \frac{(\boldsymbol{\eta}_p(\mu_k^x) \rightarrow \mu_k^x + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mu_k^x)_1}{(\boldsymbol{\eta}_p(\mu_k^x) \rightarrow \mu_k^x + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mu_k^x)_2} \\ \times \frac{(\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2)_1}{(\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2)_2} + 1 \\ -\frac{1}{2} \exp\left(\boldsymbol{\eta}_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}\right) \\ \times \frac{(\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2)_1}{(\boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2)_2} + 1 \end{bmatrix}$$

$$\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \leftarrow \left(\psi \left\{ \left(\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_k + 1 \right\} \right. \\ \left. - \psi \left[\sum_{k=1}^K \left\{ \left(\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_k + 1 \right\} \right] \right)$$

$$\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \leftarrow \exp\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}}\right)$$

$$\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} \leftarrow \alpha_k - 1 \quad ; \quad \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} \leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right) \end{bmatrix}$$

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \leftarrow \left[\begin{array}{l} \frac{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_2} \\ \times \left\{ -\frac{1}{2} \mathbf{o}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1}\right)_1}{\left(\boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1}\right)_2} \right. \\ \left. + \mathcal{J}\left(\boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}; 1, 2\right) \right\} \\ \frac{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_1^{+1}}{2\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_2} \\ \times \mathcal{K}\left(\boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1}\right) \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}$$

$$\leftarrow \left[\begin{array}{l} \frac{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_2} \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}})^\top \mathbf{o} \\ - \frac{\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_1^{+1}}{2\left(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}\right)_2} \times \text{vec} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right\} \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \leftarrow \left[-\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_\alpha^2} \mathbf{I}_2 \right) \right]$$

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \leftarrow \left[\begin{array}{l} -\frac{n}{2} \\ -\frac{1}{2} E_{q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\alpha})}^{\text{nat}} (\|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2) \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \leftarrow \left[\begin{array}{l} -\frac{3}{2} \\ \frac{\left(\boldsymbol{\eta}_{p(a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o}\right)_2} \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} \leftarrow \left[\begin{array}{l} -\frac{1}{2} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}\right)_2} \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \leftarrow \left[\begin{array}{l} -\frac{n}{2} \\ -\frac{1}{2} E_{q(\mathbf{x}_{\text{unobs}}, \boldsymbol{\beta})}^{\text{nat}} (\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2) \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} \leftarrow \left[\begin{array}{l} -\frac{1}{2} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}\right)_2} \end{array} \right]$$

$$\eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2 \leftarrow \begin{bmatrix} -\frac{1}{2}\exp\left(\eta_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}\right) \\ -\frac{1}{2}\exp\left(\eta_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}\right) \\ \times \mathcal{I}\left\{\left(\eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \eta_p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i} \right. \right. \\ \left. \left. + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow x_{\text{unobs},i}\right), \left(\eta_p(\mu_k^x) \rightarrow \mu_k^x \right. \right. \\ \left. \left. + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x\right)\right\} \end{bmatrix}$$

$$\eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik} \leftarrow -\frac{1}{2}\left[\log\left\{-\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_2\right\} \right. \\ \left. -\psi\left\{-\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_1 - 1\right\} \right. \\ \left. -\frac{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_1 + 1}{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_2}\right\} \\ \times \mathcal{I}\left\{\left(\eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \eta_p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i} \right. \right. \\ \left. \left. + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow x_{\text{unobs},i}\right), \right. \\ \left. \left(\eta_p(\mu_k^x) \rightarrow \mu_k^x + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x\right)\right\}]$$

$$\eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x \leftarrow \begin{bmatrix} -\frac{1}{2}\exp\left(\eta_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}\right) \\ \times \left\{E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{x})\right\}_i \\ \times \frac{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_1 + 1}{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_2} \\ -\frac{1}{2}\exp\left(\eta_p(a_{ik}|\omega_k) \rightarrow a_{ik} + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}\right) \\ \times \frac{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_1 + 1}{\left(\eta_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2 + \eta_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2\right)_2} \end{bmatrix}$$

$$\eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta} \leftarrow \begin{bmatrix} \frac{\left(\eta_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon + \eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2\right)_1 + 1}{\left(\eta_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon + \eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2\right)_2} \\ \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{X})^\top \mathbf{y} \\ \frac{\left(\eta_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2 + \eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2\right)_1 + 1}{2\left(\eta_p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2 + \eta_p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2\right)_2} \\ \times \text{vec}\left\{E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{X}^\top \mathbf{X})\right\} \end{bmatrix}$$

$$\begin{aligned}
 \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} &\leftarrow \left[\begin{array}{c} -\frac{3}{2} \\ \frac{\left(\boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}\right)_1 + 1}{\left(\boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}\right)_2} \end{array} \right]; \quad \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \leftarrow \left[\begin{array}{c} -\frac{3}{2} \\ -\frac{1}{A_\varepsilon^2} \end{array} \right] \\
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} &\leftarrow \left[\begin{array}{c} \frac{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}\right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}\right)_2} \\ \times \left\{ -\frac{1}{2} \mathbf{y}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x}\right)_1}{\left(\boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x}\right)_2} \right. \\ \quad \left. - \mathcal{J}\left(\boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}}; 1, 3\right) \right. \\ \quad \left. - \mathbf{c}_{x_{\text{unobs},i}} \mathcal{J}\left(\boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}}; 2, 3\right) \right\} \\ \frac{1}{4} \frac{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}\right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}\right)_2} \\ \times \frac{\left(\boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x}\right)_1}{\left(\boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x}\right)_2} \end{array} \right]
 \end{aligned}$$

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

$$\boldsymbol{\eta}_{q(\boldsymbol{\alpha})}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}^* + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}^*; \quad \boldsymbol{\eta}_{q(\sigma_o^2)}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}^* + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}^*$$

$$\boldsymbol{\eta}_{q(a_o)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o}^* + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o}^*$$

$$\begin{aligned}
 \text{For } 1 \leq k \leq K : \boldsymbol{\eta}_{q(a_k^x)}^* &\leftarrow \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x}^* + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x}^* \\
 \boldsymbol{\eta}_{q\{(\sigma_k^x)^2\}}^* &\leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2}^* + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2}^* \\
 \boldsymbol{\eta}_{q(\mu_k^x)}^* &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mu_k^x}^* + \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x}^*
 \end{aligned}$$

$$\begin{aligned}
 \text{For } 1 \leq i \leq n \text{ and } 1 \leq k \leq K : \boldsymbol{\eta}_{q(a_{ik})}^* &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}}^* + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}}^* \\
 \boldsymbol{\eta}_{q(\omega_k)}^* &\leftarrow \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k}^* + \boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k}^*
 \end{aligned}$$

$$\boldsymbol{\eta}_{q(\boldsymbol{\beta})}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}}^* + \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}^*; \quad \boldsymbol{\eta}_{q(\sigma_\varepsilon^2)}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}^* + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}^*$$

$$\boldsymbol{\eta}_{q(a_\varepsilon)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}^* + \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon}^*$$

For $1 \leq i \leq n_{\text{unobs}}$:

$$\begin{aligned}
 \boldsymbol{\eta}_{q(x_{\text{unobs},i})}^* &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}}^* + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}}^* \\
 &\quad + \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}}^*
 \end{aligned}$$

Algorithm 12: *VMP algorithm for determining the natural parameter vectors of the optimal density functions for approximate inference in Model (6.4).*

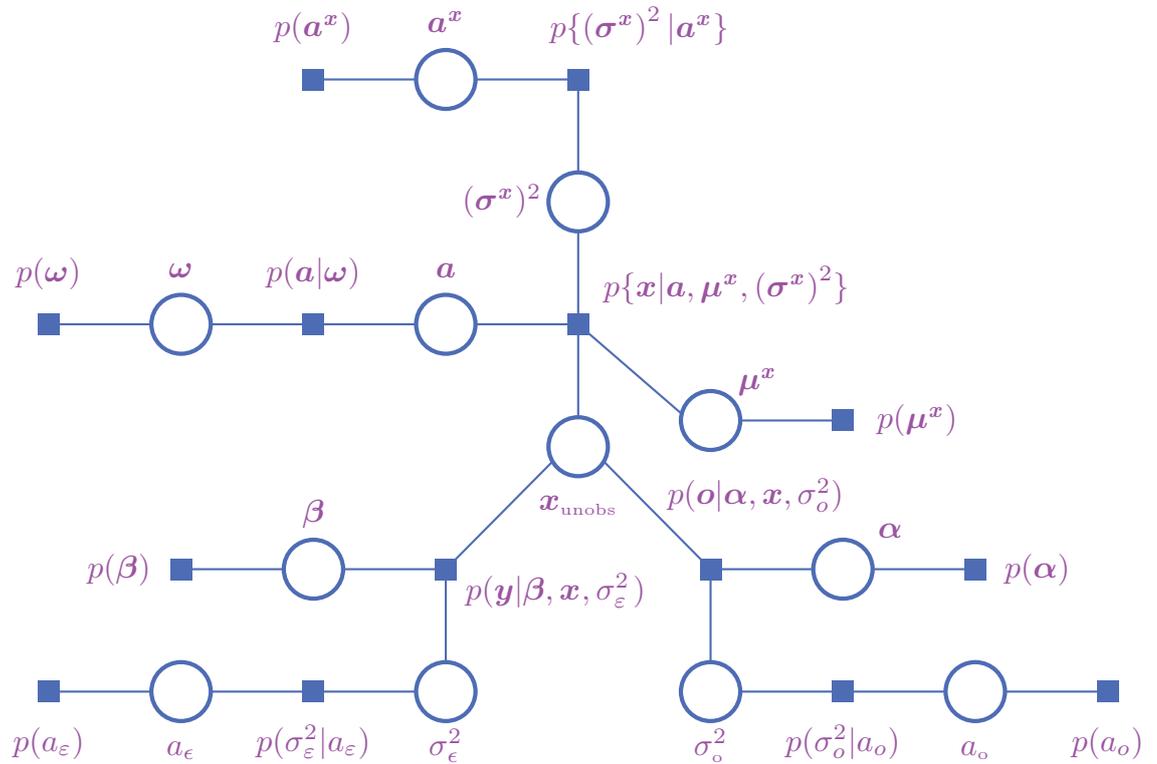


Figure 7.4: Factor graph corresponding to the mixture model with measurement error given in (6.4).

7.8 Discussion

The general VMP algorithm, presented as Algorithm 10, is a useful guide, in terms of notation, for representing the three main phases in a VMP scheme. The two examples in this chapter provided further details on the use of natural parameter forms and primitives for certain density functions when using VMP. The extension of VMP to arbitrarily large models may not be immediately obvious from the examples discussed here, however the ease by which certain primitives and results can be used for larger and more complex models makes VMP an attractive alternative to MFVB.

For the interested reader, the efficiency of the VMP approach can be seen from the examples given in Wand (2015). Wand (2015) defines the notion of a *fragment*, which is a sub-graph of a factor graph that consists only of a single factor and its corresponding stochastic neighbor nodes. These *fragments* represent the underlying structure of a wide range of large semiparametric regression models and gives VMP the advantage of being more amenable to modularization and extension to larger complex models.

7.A Derivation of Algorithm 11

In this section we describe in detail, the derivations for Algorithm 11. These derivations have made use of the natural parameter forms, primitives and functions defined in Sections 7.3, 7.4 and 7.5.

7.A.1 Step 1: Initialise factor to stochastic node messages

The derivation of the updates corresponding to the Σ node in Figure 7.3 were given in section 7.7.1. Here we give the derivations of the updates corresponding to the rest of the nodes in the factor graph. The Inverse-Gamma messages are initialised to be the Inverse-Gamma(1, 1) density function and the multivariate normal messages are initialised to be the $N(\mathbf{0}_v, \mathbf{I}_v)$ density function, where $v = p + \sum_{\ell=1}^d K_\ell$. These messages can be written as

$$\begin{aligned} \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \Sigma) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\beta}, \mathbf{u}) &\leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix} \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \Sigma) \rightarrow \boldsymbol{\nu}} \right\} \\ \mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\beta}, \mathbf{u}) &\leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix} \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} \right\} \\ \mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma^2}(\sigma^2) &\leftarrow \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma_\ell^2} \right\} \\ \mathbf{m}_{p(\sigma^2|\mathbf{a}) \rightarrow \sigma^2}(\sigma^2) &\leftarrow \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\sigma_\ell^2|\mathbf{a}_\ell) \rightarrow \sigma_\ell^2} \right\} \\ \mathbf{m}_{p(\sigma^2|\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\mathbf{a}_\ell) \\ 1/\mathbf{a}_\ell \end{bmatrix}^\top \boldsymbol{\eta}_{p(\sigma_\ell^2|\mathbf{a}_\ell) \rightarrow \mathbf{a}_\ell} \right\} \\ \mathbf{m}_{p(\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\mathbf{a}_\ell) \\ 1/\mathbf{a}_\ell \end{bmatrix}^\top \boldsymbol{\eta}_{p(\mathbf{a}_\ell) \rightarrow \mathbf{a}_\ell} \right\} \end{aligned}$$

and these settings correspond to the natural parameter vectors having the following initial values:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \Sigma) \rightarrow \boldsymbol{\nu}} \leftarrow \begin{bmatrix} \mathbf{0}_v \\ -\frac{1}{2} \text{vec}(\mathbf{I}_v) \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} \leftarrow \begin{bmatrix} \mathbf{0}_v \\ -\frac{1}{2} \text{vec}(\mathbf{I}_v) \end{bmatrix},$$

$$\begin{aligned} \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}, \\ \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}, \quad \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}. \end{aligned}$$

7.A.2 Step 2 (a): Update stochastic node to factor messages

The next stochastic node to factor message to update is $\mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})}(\boldsymbol{\nu})$, and since the only other neighbour of $\boldsymbol{\nu}$ apart from $p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})$ is $p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)$, this message takes the form

$$\mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})}(\boldsymbol{\nu}) \leftarrow \mathbf{m}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu})$$

and can be written as

$$\mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})}(\boldsymbol{\nu}) \leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix} \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} \right\}$$

where

$$\boldsymbol{\eta}_{\boldsymbol{\nu} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})} \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}}.$$

Similarly,

$$\begin{aligned} \mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)}(\boldsymbol{\nu}) &\leftarrow \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \\ \mathbf{m}_{\sigma_\ell^2 \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)}(\sigma_\ell^2) &\leftarrow \mathbf{m}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}(\sigma_\ell^2) \\ \mathbf{m}_{\sigma_\ell^2 \rightarrow p(\sigma_\ell^2|a_\ell)}(\sigma_\ell^2) &\leftarrow \mathbf{m}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2}(\sigma_\ell^2) \\ \mathbf{m}_{a_\ell \rightarrow p(\sigma_\ell^2|a_\ell)}(a_\ell) &\leftarrow \mathbf{m}_{p(a_\ell) \rightarrow a_\ell}(a_\ell) \\ \mathbf{m}_{a_\ell \rightarrow p(a_\ell)}(a_\ell) &\leftarrow \mathbf{m}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell}(a_\ell) \end{aligned}$$

which is equivalent to the following natural parameter updates:

$$\begin{aligned} \boldsymbol{\eta}_{\boldsymbol{\nu} \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} \\ \boldsymbol{\eta}_{\sigma_\ell^2 \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} \\ \boldsymbol{\eta}_{\sigma_\ell^2 \rightarrow p(\sigma_\ell^2|a_\ell)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} \\ \boldsymbol{\eta}_{a_\ell \rightarrow p(\sigma_\ell^2|a_\ell)} &\leftarrow \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} \\ \boldsymbol{\eta}_{a_\ell \rightarrow p(a_\ell)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \end{aligned}$$

7.A.3 Step 2 (b): Update factor to stochastic node messages

The next factor to stochastic node message to update is $\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu})$ and takes the form

$$\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \leftarrow \exp \left\{ \int_{\mathbb{R}^{n \times n}} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\Sigma} \rightarrow p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma})}(\boldsymbol{\Sigma}) \times \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) \times \log p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) d\boldsymbol{\Sigma} \right\}$$

where Z is the normalising constant.

$$\begin{aligned} & \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \\ & \propto \exp \left[\int_{\mathbb{R}^{n \times n}} \left\{ \left[\begin{array}{c} \log |\boldsymbol{\Sigma}| \\ \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{array} \right]^{\top} \left(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} \right) \right\} \right. \\ & \quad \left. \times \left\{ \left[\begin{array}{c} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^{\top}) \end{array} \right]^{\top} \left[\begin{array}{c} \mathbf{C}^{\top} (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{y} \\ \frac{1}{2} \text{vec}(\mathbf{C}^{\top} (\mathbf{I}_m \otimes \boldsymbol{\Sigma}^{-1}) \mathbf{C}) \end{array} \right] \right\} d\boldsymbol{\Sigma} \right] \\ & = \exp \left\{ \left[\begin{array}{c} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^{\top}) \end{array} \right]^{\top} \left[\begin{array}{c} \mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \int_{\mathbb{R}^{n \times n}} \boldsymbol{\Sigma}^{-1} p_{\text{IW}_{\text{nat}}}(\boldsymbol{\Sigma}; \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) d\boldsymbol{\Sigma} \right\} \mathbf{y} \\ \frac{1}{2} \text{vec} \left[\mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \int_{\mathbb{R}^{n \times n}} \boldsymbol{\Sigma}^{-1} p_{\text{IW}_{\text{nat}}}(\boldsymbol{\Sigma}; \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) d\boldsymbol{\Sigma} \right\} \mathbf{C} \right] \end{array} \right] \right\} \\ & = \exp \left\{ \left[\begin{array}{c} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^{\top}) \end{array} \right]^{\top} \left[\begin{array}{c} \mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) \right\} \mathbf{y} \\ \frac{1}{2} \text{vec} \left[\mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) \right\} \mathbf{C} \right] \end{array} \right] \right\}. \end{aligned}$$

This means that the natural parameter vector has the following update:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} \leftarrow \left[\begin{array}{c} \mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) \right\} \mathbf{y} \\ \frac{1}{2} \text{vec} \left[\mathbf{C}^{\top} \left\{ \mathbf{I}_m \otimes \mathcal{G}(\boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}) \right\} \mathbf{C} \right] \end{array} \right].$$

The next message begins with

$$\mathbf{m}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}^d} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\sigma}^2 \rightarrow p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2)}(\boldsymbol{\sigma}^2) \times \mathbf{m}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\sigma}^2}(\boldsymbol{\sigma}^2) \times \log p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) d\boldsymbol{\sigma}^2 \right\}$$

Collecting like terms in $\mathbf{m}_{\sigma^2 \rightarrow p(\boldsymbol{\nu}|\sigma^2)}(\sigma^2)$ and $\mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma^2}(\sigma^2)$ gives:

$$\begin{aligned}
 & \mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \\
 & \propto \exp \left[\int_{\mathbb{R}_{\geq 0}^d} \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma_\ell^2} \right) \right\} \right. \\
 & \quad \left. \times \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix}^\top \left[-\frac{1}{2} \text{vec} \left(\begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left(\frac{1}{\sigma_\ell^2} \mathbf{I}_{K_\ell} \right) \end{bmatrix} \right) \right] \right\} d\boldsymbol{\Sigma} \right] \\
 & = \exp \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix}^\top \right. \\
 & \quad \left. \times \left[-\frac{1}{2} \text{vec} \left(\begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left\{ \left(\int_0^\infty \frac{1}{\sigma_\ell^2} p_{\text{IG}_{\text{nat}}}(\boldsymbol{\eta}^\oplus) d\sigma_\ell^2 \right) \mathbf{I}_{K_\ell} \right\} \end{bmatrix} \right) \right] \right\}
 \end{aligned}$$

where $\boldsymbol{\eta}^\oplus \equiv \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma_\ell^2} \equiv [\boldsymbol{\eta}_1^\oplus \ \boldsymbol{\eta}_2^\oplus]^\top$. Using Primitive 7.4.6, we get

$$= \exp \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix}^\top \left[-\frac{1}{2} \text{vec} \left(\begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left\{ \left(\frac{\boldsymbol{\eta}_1^\oplus + 1}{\boldsymbol{\eta}_2^\oplus} \right) \mathbf{I}_{K_\ell} \right\} \end{bmatrix} \right) \right] \right\}.$$

This leads to the natural parameter update

$$\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} \leftarrow \left[-\frac{1}{2} \text{vec} \left(\begin{bmatrix} \mathbf{F}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{blockdiag} \left\{ \left(\frac{\boldsymbol{\eta}_1^\oplus + 1}{\boldsymbol{\eta}_2^\oplus} \right) \mathbf{I}_{K_\ell} \right\} \end{bmatrix} \right) \right].$$

Next we look at the factor to stochastic node message going from $p(\boldsymbol{\nu}|\sigma^2)$ to σ^2 :

$$\begin{aligned}
 \mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \sigma^2}(\sigma^2) & \leftarrow \exp \left\{ \int_{\mathbb{R}^v} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\nu} \rightarrow p(\boldsymbol{\nu}|\sigma^2)}(\boldsymbol{\nu}) \times \mathbf{m}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\nu}) \right. \\
 & \quad \left. \times \log p(\boldsymbol{\nu}|\sigma^2) d\boldsymbol{\nu} \right\} \\
 & \propto \exp \left[\int_{\mathbb{R}^v} \left\{ \begin{bmatrix} \boldsymbol{\nu} \\ \text{vec}(\boldsymbol{\nu}\boldsymbol{\nu}^\top) \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\sigma^2) \rightarrow \boldsymbol{\nu}} \right) \right\} \right. \\
 & \quad \left. \times \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2} K_\ell \\ -\frac{1}{2} \|\mathbf{u}_\ell\|^2 \end{bmatrix} \right\} d\boldsymbol{\nu} \right].
 \end{aligned}$$

7.A. DERIVATION OF ALGORITHM 11

Since $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}}$ is a vector that can be partitioned by each of its components, i.e.,

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\nu}} \equiv \begin{bmatrix} \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \beta} \\ \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_1} \\ \vdots \\ \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_d} \end{bmatrix},$$

and similarly for $\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}}$, then we can write $\mathbf{m}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}}(\boldsymbol{\sigma}^2)$ as

$$\exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \times \left[-\frac{1}{2} \int_{\mathbb{R}^{\mathbf{u}_\ell}} \|\mathbf{u}_\ell\|^2 p_{\text{Nnat,vec}} \left(\mathbf{u}_\ell; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_\ell} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \mathbf{u}_\ell} \right) d\mathbf{u}_\ell \right] \right\}.$$

Using Primitive 7.4.3, this becomes

$$\exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \left[-\frac{1}{2} \mathcal{F} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_\ell} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \mathbf{u}_\ell} \right) \right] \right\}$$

and thus the natural parameter update is

$$\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \boldsymbol{\nu}} \leftarrow \left[-\frac{1}{2} \mathcal{F} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \mathbf{u}_\ell} + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \mathbf{u}_\ell} \right) \right].$$

The next message to update is

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\sigma}^2|\mathbf{a}) \rightarrow \boldsymbol{\sigma}^2}(\boldsymbol{\sigma}^2) &\leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}^d} \frac{1}{Z} \mathbf{m}_{\mathbf{a} \rightarrow p(\boldsymbol{\sigma}^2|\mathbf{a})}(\mathbf{a}) \times \mathbf{m}_{p(\boldsymbol{\sigma}^2|\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a}) \right. \\ &\quad \left. \times \log p(\boldsymbol{\sigma}^2|\mathbf{a}) d\mathbf{a} \right\} \\ &\propto \exp \left\{ \int_{\mathbb{R}_{\geq 0}^d} \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(a_\ell) \\ 1/a_\ell \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \right) \right\} \right. \\ &\quad \left. \times \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -1/a_\ell \end{bmatrix} \right\} d\mathbf{a} \right\} \\ &= \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \left[-\int_0^\infty \frac{1}{a_\ell} p_{\text{IGnat}} \left(a_\ell; \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \right) da_\ell \right] \right\}. \end{aligned}$$

The use of Primitive 7.4.6 gives

$$\mathbf{m}_{p(\boldsymbol{\sigma}^2|\mathbf{a}) \rightarrow \boldsymbol{\sigma}^2}(\boldsymbol{\sigma}^2) \propto \exp \left\{ \sum_{\ell=1}^d \begin{bmatrix} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{bmatrix}^\top \left[\frac{-3/2}{\left(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \right)_1} + 1 \right] \right\}.$$

Therefore, the natural parameter update is given as

$$\boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} \leftarrow \left[\begin{array}{c} -3/2 \\ -\frac{\left(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell}\right)_2} \end{array} \right].$$

The second last factor to stochastic node message to update is

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\sigma}^2|\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}^d} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\sigma}^2 \rightarrow p(\boldsymbol{\sigma}^2|\mathbf{a})}(\boldsymbol{\sigma}^2) \times \mathbf{m}_{p(\boldsymbol{\sigma}^2|\mathbf{a}) \rightarrow \boldsymbol{\sigma}^2}(\boldsymbol{\sigma}^2) \right. \\ &\quad \left. \times \log p(\boldsymbol{\sigma}^2|\mathbf{a}) d\boldsymbol{\sigma}^2 \right\} \\ &\propto \exp \left\{ \int_{\mathbb{R}_{\geq 0}^d} \left\{ \sum_{\ell=1}^d \left[\begin{array}{c} \log(\sigma_\ell^2) \\ 1/\sigma_\ell^2 \end{array} \right]^\top \left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2} \right) \right\} \right. \\ &\quad \left. \times \left\{ \sum_{\ell=1}^d \left[\begin{array}{c} \log(a_\ell) \\ 1/a_\ell \end{array} \right]^\top \left[\begin{array}{c} -\frac{1}{2} \\ -1/\sigma_\ell^2 \end{array} \right] \right\} d\sigma_\ell^2 \right\} \\ &= \exp \left\{ \sum_{\ell=1}^d \left[\begin{array}{c} \log(a_\ell) \\ 1/a_\ell \end{array} \right]^\top \left[-\int_0^\infty \frac{1}{\sigma_\ell^2} p_{\text{IG}_{\text{nat}}}(\sigma_\ell^2; \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}) d\sigma_\ell^2 \right] \right\} \end{aligned}$$

Using Primitive 7.4.6, this becomes

$$\exp \left\{ \sum_{\ell=1}^d \left[\begin{array}{c} \log(a_\ell) \\ 1/a_\ell \end{array} \right]^\top \left[\begin{array}{c} -\frac{1}{2} \\ -\frac{\left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}\right)_2} \end{array} \right] \right\},$$

and so the natural parameter update is

$$\boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell} \leftarrow \left[\begin{array}{c} -\frac{1}{2} \\ -\frac{\left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}\right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2} + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}\right)_2} \end{array} \right].$$

The last factor to stochastic node message to update is $\mathbf{m}_{p(\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a})$ and has the form

$$\begin{aligned} \mathbf{m}_{p(\mathbf{a}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \frac{1}{Z} \log p(\mathbf{a}) \right\} \\ &\propto \exp \left\{ \sum_{\ell=1}^d \left[\begin{array}{c} \log(a_\ell) \\ 1/a_\ell \end{array} \right]^\top \left[\begin{array}{c} -3/2 \\ -1/A_\ell^2 \end{array} \right] \right\} \end{aligned}$$

which means that the natural parameter update corresponding to this message is

$$\boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell} \leftarrow \left[\begin{array}{c} -3/2 \\ -1/A_\ell^2 \end{array} \right]$$

and remains fixed at this constant value since A_ℓ is a hyperparameter to be specified by the user.

7.A.4 Step 3 : Optimal q -densities

Once convergence is achieved from Steps 2(a) and (b) we get the optimal q -densities of each stochastic node through:

$$\boldsymbol{\eta}_{q(\boldsymbol{\nu})}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow (\boldsymbol{\nu})}^* + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow (\boldsymbol{\nu})}^*$$

$$\boldsymbol{\eta}_{q(\boldsymbol{\Sigma})}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\nu}, \boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}^* + \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}) \rightarrow \boldsymbol{\Sigma}}^*$$

and for $1 \leq \ell \leq d$:

$$\boldsymbol{\eta}_{q(a_\ell)}^* \leftarrow \boldsymbol{\eta}_{p(a_\ell) \rightarrow a_\ell}^* + \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow a_\ell}^*$$

$$\boldsymbol{\eta}_{q(\sigma_\ell^2)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_\ell^2|a_\ell) \rightarrow \sigma_\ell^2}^* + \boldsymbol{\eta}_{p(\boldsymbol{\nu}|\boldsymbol{\sigma}^2) \rightarrow \sigma_\ell^2}^*.$$

7.B Derivation of Algorithm 12

This section gives the derivations for the updates given in Algorithm 12. Here, we also make use of the natural parameter forms, primitives, results and functions defined in Sections 7.3, 7.4 and 7.5.

7.B.1 Step 1: Initialise factor to stochastic node messages

We initialise the factor to stochastic node messages to be

$$\begin{aligned}
\mathbf{m}_{p(\mathbf{a}^x) \rightarrow \mathbf{a}^x}(\mathbf{a}^x) &\leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log(a_k^x) \\ 1/a_k^x \end{bmatrix}^\top \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \right\} \\
\mathbf{m}_{p\{(\sigma^x)^2 | \mathbf{a}^x\} \rightarrow \mathbf{a}^x}(\mathbf{a}^x) &\leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log(a_k^x) \\ 1/a_k^x \end{bmatrix}^\top \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} \right\} \\
\mathbf{m}_{p\{(\sigma^x)^2 | \mathbf{a}^x\} \rightarrow (\sigma^x)^2} \{(\sigma^x)^2\} &\leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log\{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow (\sigma_k^x)^2} \right\} \\
\mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma^x)^2} \{(\sigma^x)^2\} &\leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log\{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix}^\top \right. \\
&\quad \left. \times \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2} \right\} \\
\mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \sum_{i=1}^n \begin{bmatrix} a_{i1} \\ \vdots \\ a_{iK} \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{i1}} \\ \vdots \\ \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{iK}} \end{bmatrix} \right\} \\
\mathbf{m}_{p(\mathbf{a} | \boldsymbol{\omega}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \sum_{i=1}^n \begin{bmatrix} a_{i1} \\ \vdots \\ a_{iK} \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\eta}_{p(a_{i1} | \omega_1) \rightarrow a_{i1}} \\ \vdots \\ \boldsymbol{\eta}_{p(a_{iK} | \omega_K) \rightarrow a_{iK}} \end{bmatrix} \right\} \\
\mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) &\leftarrow \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \right. \\
&\quad \left. \times \boldsymbol{\eta}_{p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}} \right\} \\
\mathbf{m}_{p(\boldsymbol{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) &\leftarrow \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \right\}
\end{aligned}$$

$$\mathbf{m}_{p\{\mathbf{x}|\boldsymbol{\alpha}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \boldsymbol{\mu}^x}(\boldsymbol{\mu}^x) \leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \mu_k \\ \mu_k^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p\{\mathbf{x}|\boldsymbol{\alpha}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{\mu}^x) \rightarrow \boldsymbol{\mu}^x}(\boldsymbol{\mu}^x) \leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \mu_k \\ \mu_k^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \right\}$$

$$\mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) \leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log(\omega_1) \\ \vdots \\ \log(\omega_K) \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\eta}_{p(\mathbf{a}_1|\omega_1) \rightarrow \omega_1} \\ \vdots \\ \boldsymbol{\eta}_{p(\mathbf{a}_K|\omega_K) \rightarrow \omega_K} \end{bmatrix} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) \leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log(\omega_1) \\ \vdots \\ \log(\omega_K) \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\eta}_{p(\omega_1) \rightarrow \omega_1} \\ \vdots \\ \boldsymbol{\eta}_{p(\omega_K) \rightarrow \omega_K} \end{bmatrix} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) \leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) \leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2) \leftarrow \exp \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \right\}$$

$$\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}(\sigma_o^2) \leftarrow \exp \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right\}$$

$$\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o}(a_o) \leftarrow \exp \left\{ \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} \right\}$$

$$\mathbf{m}_{p(a_o) \rightarrow a_o}(a_o) \leftarrow \exp \left\{ \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \boldsymbol{\eta}_{p(a_o) \rightarrow a_o} \right\}$$

$$\mathbf{m}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp \left\{ \begin{bmatrix} \boldsymbol{\beta} \\ \text{vec}(\boldsymbol{\beta}\boldsymbol{\beta}^\top) \end{bmatrix}^\top \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} \right\}$$

$$\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \leftarrow \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} \right\}$$

$$\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \leftarrow \exp \left\{ \left[\begin{array}{c} \boldsymbol{\beta} \\ \text{vec}(\boldsymbol{\beta}\boldsymbol{\beta}^\top) \end{array} \right]^\top \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} \right\}$$

$$\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2) \leftarrow \exp \left\{ \left[\begin{array}{c} \log(\sigma_\varepsilon^2) \\ 1/\sigma_\varepsilon^2 \end{array} \right]^\top \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \right\}$$

$$\mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2) \leftarrow \exp \left\{ \left[\begin{array}{c} \log(\sigma_\varepsilon^2) \\ 1/\sigma_\varepsilon^2 \end{array} \right]^\top \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right\}$$

$$\mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon) \leftarrow \exp \left\{ \left[\begin{array}{c} \log(a_\varepsilon) \\ 1/a_\varepsilon \end{array} \right]^\top \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} \right\}$$

$$\mathbf{m}_{p(a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon) \leftarrow \exp \left\{ \left[\begin{array}{c} \log(a_\varepsilon) \\ 1/a_\varepsilon \end{array} \right]^\top \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \right\}.$$

The initial values for the corresponding natural parameter vectors are set as:

$$\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}$$

$$\boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix};$$

$$\boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}} \leftarrow \log\left(\frac{1}{K}\right); \quad \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \leftarrow \log\left(\frac{1}{K}\right);$$

$$\boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mu_k^x} \leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; \quad \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; \quad \boldsymbol{\eta}_{p(\mathbf{a}_k|\omega_k) \rightarrow \omega_k} \leftarrow 1;$$

$$\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} \leftarrow \alpha_k - 1; \quad \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}} \leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix};$$

$$\boldsymbol{\eta}_{p(o|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; \quad \boldsymbol{\eta}_{p(o|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \leftarrow \begin{bmatrix} \mathbf{0}_2 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_2) \end{bmatrix};$$

$$\boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \leftarrow \begin{bmatrix} \mathbf{0}_2 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_2) \end{bmatrix}; \quad \boldsymbol{\eta}_{p(o|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix};$$

$$\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\eta}_{p(a_o) \rightarrow a_o} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; \quad \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} \leftarrow \begin{bmatrix} \mathbf{0}_3 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_3) \end{bmatrix}$$

$$\begin{aligned}
 \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow x_{\text{unobs},i} &\leftarrow \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}; & \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \boldsymbol{\beta} &\leftarrow \begin{bmatrix} \mathbf{0}_3 \\ -\frac{1}{2}\text{vec}(\mathbf{I}_3) \end{bmatrix} \\
 \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow a_\varepsilon &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}; & \boldsymbol{\eta}_{p(a_\varepsilon)} \rightarrow a_\varepsilon &\leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}
 \end{aligned}$$

7.B.2 Step 2 (a): Update stochastic node to factor messages

The stochastic node to factor messages have the following updates:

$$\begin{aligned}
 \mathbf{m}_{\mathbf{a}^x \rightarrow p(\mathbf{a}^x)}(\mathbf{a}^x) &\leftarrow \mathbf{m}_{p\{(\sigma^x)^2|\mathbf{a}^x\} \rightarrow \mathbf{a}^x}(\mathbf{a}^x) \\
 \mathbf{m}_{(\sigma^x)^2 \rightarrow p\{(\sigma^x)^2|\mathbf{a}^x\}} &\leftarrow \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma^x)^2} \\
 \mathbf{m}_{(\sigma^x)^2 \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}} &\leftarrow \mathbf{m}_{p\{(\sigma^x)^2|\mathbf{a}^x\} \rightarrow (\sigma^x)^2} \\
 \mathbf{m}_{\mathbf{a} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}} &\leftarrow \mathbf{m}_{p(\mathbf{a}) \rightarrow \mathbf{a}}; & \mathbf{m}_{\mathbf{a}} \rightarrow p(\mathbf{a}) &\leftarrow \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{a}} \\
 \mathbf{m}_{\boldsymbol{\omega} \rightarrow p(\mathbf{a}|\boldsymbol{\omega})} &\leftarrow \mathbf{m}_{p(\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}; & \mathbf{m}_{\boldsymbol{\omega}} \rightarrow p(\boldsymbol{\omega}) &\leftarrow \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}} \\
 \mathbf{m}_{\boldsymbol{\mu}^x \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}} &\leftarrow \mathbf{m}_{p(\boldsymbol{\mu}^x) \rightarrow \boldsymbol{\mu}^x} \\
 \mathbf{m}_{\boldsymbol{\mu}^x \rightarrow p(\boldsymbol{\mu}^x)} &\leftarrow \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \boldsymbol{\mu}^x} \\
 \mathbf{m}_{\boldsymbol{\alpha} \rightarrow p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \mathbf{m}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}; & \mathbf{m}_{\boldsymbol{\alpha}} \rightarrow p(\boldsymbol{\alpha}) &\leftarrow \mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \\
 \mathbf{m}_{\sigma_o^2 \rightarrow p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}; & \mathbf{m}_{\sigma_o^2} \rightarrow p(\sigma_o^2|a_o) &\leftarrow \mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \\
 \mathbf{m}_{a_o \rightarrow p(\sigma_o^2|a_o)} &\leftarrow \mathbf{m}_{p(a_o) \rightarrow a_o}; & \mathbf{m}_{a_o} \rightarrow p(a_o) &\leftarrow \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o} \\
 \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}} &\leftarrow \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}}} + \mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}} \\
 \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} &\leftarrow \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}}} + \mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}} \\
 \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}}} + \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}}} \\
 \mathbf{m}_{\boldsymbol{\beta} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} &\leftarrow \mathbf{m}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}; & \mathbf{m}_{\boldsymbol{\beta}} \rightarrow p(\boldsymbol{\beta}) &\leftarrow \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} \\
 \mathbf{m}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} &\leftarrow \mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}; & \mathbf{m}_{\sigma_\varepsilon^2} \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon) &\leftarrow \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \\
 \mathbf{m}_{a_\varepsilon \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} &\leftarrow \mathbf{m}_{p(a_\varepsilon) \rightarrow a_\varepsilon}; & \mathbf{m}_{a_\varepsilon} \rightarrow p(a_\varepsilon) &\leftarrow \mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}
 \end{aligned}$$

and this is equivalent to the following natural parameter updates:

$$\begin{aligned}
 \boldsymbol{\eta}_{a_k^x \rightarrow p(a_k^x)} &\leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow a_k^x}; & \boldsymbol{\eta}_{a_k^x} \rightarrow p\{(\sigma_k^x)^2|a_k^x\} &\leftarrow \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \\
 \boldsymbol{\eta}_{(\sigma_k^x)^2 \rightarrow p\{(\sigma_k^x)^2|a_k^x\}} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow (\sigma_k^x)^2} \\
 \boldsymbol{\eta}_{(\sigma_k^x)^2 \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma_k^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2} \\
 \boldsymbol{\eta}_{a_{ik}} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma_k^x)^2\} &\leftarrow \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}}; & \boldsymbol{\eta}_{a_{ik}} \rightarrow p(a_{ik}|\omega_k) &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}} \\
 \boldsymbol{\eta}_{\omega_k \rightarrow p(a_{ik}|\omega_k)} &\leftarrow \boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k}; & \boldsymbol{\eta}_{\omega_k} \rightarrow p(\omega_k) &\leftarrow \boldsymbol{\eta}_{p(\mathbf{a}_k|\omega_k) \rightarrow \omega_k}
 \end{aligned}$$

$$\begin{aligned}
 \boldsymbol{\eta}_{\mu_k^x \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma_k^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x}; & \boldsymbol{\eta}_{\mu_k^x \rightarrow p(\mu_k^x)} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mu_k^x} \\
 \boldsymbol{\eta}_{\boldsymbol{\alpha} \rightarrow p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}; & \boldsymbol{\eta}_{\boldsymbol{\alpha} \rightarrow p(\boldsymbol{\alpha})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \\
 \boldsymbol{\eta}_{\sigma_o^2 \rightarrow p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}; & \boldsymbol{\eta}_{\sigma_o^2 \rightarrow p(\sigma_o^2|a_o)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} \\
 \boldsymbol{\eta}_{a_o \rightarrow p(\sigma_o^2|a_o)} &\leftarrow \boldsymbol{\eta}_{p(a_o) \rightarrow a_o}; & \boldsymbol{\eta}_{a_o \rightarrow p(a_o)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} \\
 \boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \\
 \boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \\
 \boldsymbol{\eta}_{x_{\text{unobs},i} \rightarrow p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)} &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} \\
 \boldsymbol{\eta}_{\boldsymbol{\beta} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} &\leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}; & \boldsymbol{\eta}_{\boldsymbol{\beta} \rightarrow p(\boldsymbol{\beta})} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} \\
 \boldsymbol{\eta}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \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}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} \\
 \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(\sigma_\varepsilon^2|a_\varepsilon)} &\leftarrow \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon}; & \boldsymbol{\eta}_{a_\varepsilon \rightarrow p(a_\varepsilon)} &\leftarrow \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}
 \end{aligned}$$

7.B.3 Step 2 (b): Update factor to stochastic node messages

We begin with the update for the factor to stochastic node message $\mathbf{m}_{p(\mathbf{a}^x) \rightarrow \mathbf{a}^x}(\mathbf{a}^x)$. Since the only neighbour of factor $p(\mathbf{a}^x)$ is the stochastic node \mathbf{a}^x this message takes on the form

$$\begin{aligned}
 \mathbf{m}_{p(\mathbf{a}^x) \rightarrow \mathbf{a}^x}(\mathbf{a}^x) &\leftarrow \exp\left\{\frac{1}{Z} \log p(\mathbf{a}^x)\right\} \\
 &\propto \exp\left\{\sum_{k=1}^K \begin{bmatrix} \log(a_k^x) \\ 1/a_k^x \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -1/(A_k^x)^2 \end{bmatrix}\right\}.
 \end{aligned}$$

This means that the update for the natural parameter vector is

$$\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \leftarrow \begin{bmatrix} -3/2 \\ -1/(A_k^x)^2 \end{bmatrix}$$

and remains constant at this value since A_k^x is a hyperparameter to be specified by the user. Next we consider the message going from factor $p\{(\sigma^x)^2|\mathbf{a}^x\}$ to stochastic node \mathbf{a}^x :

$$\begin{aligned}
 \mathbf{m}_{p\{(\sigma^x)^2|\mathbf{a}^x\} \rightarrow \mathbf{a}^x}(\mathbf{a}^x) &\leftarrow \exp\left[\int_{\mathbb{R}_{>0}^K} \frac{1}{Z} \mathbf{m}_{(\sigma^x)^2 \rightarrow p\{(\sigma^x)^2|\mathbf{a}^x\}}\{(\sigma^x)^2\} \right. \\
 &\quad \times \mathbf{m}_{p\{(\sigma^x)^2|\mathbf{a}^x\} \rightarrow (\sigma^x)^2}\{(\sigma^x)^2\} \\
 &\quad \left. \times \log p\{(\sigma^x)^2|\mathbf{a}^x\} d(\sigma^2)^2\right]
 \end{aligned}$$

$$\begin{aligned}
 &\propto \exp \left[\int_{\mathbb{R}_{\geq 0}^K} \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix} \right\}^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right. \\
 &\quad \times \left. \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{a_k^x\} \\ 1/a_k^x \end{bmatrix} \right\}^\top \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{(\sigma_k^x)^2} \end{bmatrix} \right] d(\boldsymbol{\sigma}^x)^2 \\
 &\propto \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{a_k^x\} \\ 1/a_k^x \end{bmatrix} \right\}^\top \left[-\int_0^\infty \frac{1}{(\sigma_k^x)^2} p_{\text{IGnat}} \left((\sigma_k^x)^2 ; \boldsymbol{\eta}^\boxplus \right) d(\sigma_k^x)^2 \right] \right\}
 \end{aligned}$$

where $\boldsymbol{\eta}^\boxplus \equiv \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \equiv [\eta_1^\boxplus \ \eta_2^\boxplus]^\top$. Using Primitive 7.4.6, we get

$$\mathbf{m}_{p \{ (\sigma^x)^2 | a^x \} \rightarrow a^x (a^x)} \leftarrow \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{a_k^x\} \\ 1/a_k^x \end{bmatrix} \right\}^\top \begin{bmatrix} -\frac{1}{2} \\ -\frac{\eta_1^\boxplus + 1}{\eta_2^\boxplus} \end{bmatrix} \right\}$$

and thus, the natural parameter update is

$$\boldsymbol{\eta}_{p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow a_k^x} \leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\frac{\eta_1^\boxplus + 1}{\eta_2^\boxplus} \end{bmatrix}.$$

The next message is expressed as

$$\begin{aligned}
 \mathbf{m}_{p \{ (\sigma^x)^2 | a^x \} \rightarrow (\boldsymbol{\sigma}^x)^2 \{ (\boldsymbol{\sigma}^x)^2 \}} &\leftarrow \exp \left[\int_{\mathbb{R}_{\geq 0}^K} \frac{1}{Z} \mathbf{m}_{a^x \rightarrow p \{ (\sigma^x)^2 | a^x \}} (a^x) \right. \\
 &\quad \times \mathbf{m}_{p \{ (\sigma^x)^2 | a^x \} \rightarrow a^x (a^x)} \\
 &\quad \times \log p \{ (\boldsymbol{\sigma}^x)^2 | a^x \} d a^x \\
 &\left. \propto \exp \left[\int_{\mathbb{R}_{\geq 0}^K} \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{a_k^x\} \\ 1/a_k^x \end{bmatrix} \right\}^\top \left(\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow a_k^x \right) \right. \right. \\
 &\quad \times \left. \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix} \right\}^\top \begin{bmatrix} -3/2 \\ -\frac{1}{a_k^x} \end{bmatrix} \right] d a^x \\
 &\left. \propto \exp \left\{ \sum_{k=1}^K \begin{bmatrix} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix} \right\}^\top \left[-\int_0^\infty \frac{1}{a_k^x} p_{\text{IGnat}} \left(a_k^x ; \boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} \right) \right. \right. \\
 &\quad \left. \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow a_k^x \right] d a_k^x \right\}.
 \end{aligned}$$

Using Primitive 7.4.6, this gives

$$\begin{aligned} & \mathbf{m}_{p\{(\sigma^x)^2 | \mathbf{a}^x\} \rightarrow (\sigma^x)^2 \{(\sigma^x)^2\}} \\ & \leftarrow \exp \left\{ \sum_{k=1}^K \left[\begin{array}{c} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{array} \right]^\top \left[\begin{array}{c} -3/2 \\ \left(\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} \right)_1 + 1 \\ \left(\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} \right)_2 \end{array} \right] \right\} \end{aligned}$$

and so, the natural parameter update is

$$\boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow (\sigma_k^x)^2} \leftarrow \left[\begin{array}{c} -3/2 \\ \left(\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} \right)_1 + 1 \\ \left(\boldsymbol{\eta}_{p(a_k^x) \rightarrow a_k^x} + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2 | a_k^x\} \rightarrow a_k^x} \right)_2 \end{array} \right].$$

Next,

$$\begin{aligned} \mathbf{m}_{p(\mu^x) \rightarrow \mu^x (\mu^x)} & \leftarrow \exp \left\{ \frac{1}{Z} \log p(\mu^x) \right\} \\ & \propto \exp \left\{ \sum_{k=1}^K \left[\begin{array}{c} \mu_k^x \\ (\mu_k^x)^2 \end{array} \right]^\top \left[\begin{array}{c} \mu_\mu \sigma_\mu^2 \\ -1/2\sigma_\mu^2 \end{array} \right] \right\}. \end{aligned}$$

So, the constant value for the natural parameter vector is

$$\boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \leftarrow \left[\begin{array}{c} \mu_\mu \sigma_\mu^2 \\ -1/2\sigma_\mu^2 \end{array} \right].$$

Moving along, we next consider the message going from factor $p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\}$ to the stochastic node $\mathbf{x}_{\text{unobs}}$:

$$\begin{aligned} & \mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}} (\mathbf{x}_{\text{unobs}})} \\ & \leftarrow \exp \left[\sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \frac{1}{Z} \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} \mathbf{m}_{\mathbf{a} \rightarrow p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\}}(\mathbf{a}) \right. \\ & \quad \times \mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\} \rightarrow \mathbf{a}}(\mathbf{a}) \times \mathbf{m}_{\mu^x \rightarrow p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\}}(\mu^x) \\ & \quad \times \mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\} \rightarrow \mu^x (\mu^x)} \times \mathbf{m}_{(\sigma^x)^2 \rightarrow p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\}}\{(\sigma^x)^2\} \\ & \quad \times \mathbf{m}_{p\{\mathbf{x} | \mathbf{a}, \mu^x, (\sigma^x)^2\} \rightarrow (\sigma^x)^2 \{(\sigma^x)^2\}} \times \log p\{\mathbf{a} | \mathbf{a}, \mu^x, (\sigma^x)^2\} \\ & \quad \left. d\mathbf{a} d\mu^x d(\sigma^x)^2 \right] \end{aligned}$$

$$\begin{aligned}
 & \propto \exp \left[\sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} \left\{ a_{ik} \left(\eta_{p(a_{ik}|\omega_k)} \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \right\} \right. \\
 & \quad \times \left\{ \left[\begin{array}{c} \mu_k^x \\ (\mu_k^x)^2 \end{array} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \boldsymbol{\eta}_p (\mu_k^x) \rightarrow \mu_k^x \right) \right\} \\
 & \quad \times \left\{ \left[\begin{array}{c} \log \{ (\sigma_k^x)^2 \} \\ 1/(\sigma_k^x)^2 \end{array} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right\} \\
 & \quad \times \left[\begin{array}{c} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{array} \right]^\top \left[\begin{array}{c} a_{ik} \mu_k^x / (\sigma_k^x)^2 \\ -a_{ik} / 2 (\sigma_k^x)^2 \end{array} \right] d \mathbf{a} d \boldsymbol{\mu}^x d (\boldsymbol{\sigma}^x)^2 \Big] \\
 & = \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \sum_{k=1}^K \left[\begin{array}{c} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{array} \right]^\top \right. \\
 & \quad \times \left. \left[\begin{array}{l} \sum_{a_{ik}=0}^1 a_{ik} \times a_{ik} \times \left(\eta_{p(a_{ik}|\omega_k)} \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \right. \\ \quad \times \int_{\mathbb{R}} \mu_k^x p_{\text{Nnat}} \left(\mu_k^x ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x \right. \\ \quad \left. \left. + \boldsymbol{\eta}_p (\mu_k^x) \rightarrow \mu_k^x \right) d \mu_k^x \times \int_{\mathbb{R}_{\geq 0}} \frac{1}{(\sigma_k^x)^2} p_{\text{IGnat}} \left((\sigma_k^x)^2 ; \right. \right. \\ \quad \left. \left. \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right. \\ \quad \left. d (\sigma_k^x)^2 \right. \\ \\ \quad \left. - \frac{1}{2} \sum_{a_{ik}=0}^1 a_{ik} \times a_{ik} \times \left(\eta_{p(a_{ik}|\omega_k)} \rightarrow a_{ik} \right. \right. \\ \quad \left. \left. + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \times \int_{\mathbb{R}_{\geq 0}} \frac{1}{(\sigma_k^x)^2} p_{\text{IGnat}} \left((\sigma_k^x)^2 ; \right. \right. \\ \quad \left. \left. \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right. \\ \quad \left. d (\sigma_k^x)^2 \right] \Big] .
 \end{aligned}$$

This gives

7.B. DERIVATION OF ALGORITHM 12

$$\begin{aligned}
 \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{x}_{\text{unobs}} (\mathbf{x}_{\text{unobs}}) &\leftarrow \\
 = \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \sum_{k=1}^K \left[\begin{array}{c} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{array} \right]^\top \right. & \\
 \times \left[\begin{array}{l} -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right)_1}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right)_2} \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1^{+1}}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \\ -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1^{+1}}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \end{array} \right] & \\
 \left. \right\} &
 \end{aligned}$$

and so the update for the corresponding natural parameter is

$$\begin{aligned}
 \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} &\leftarrow \\
 \left[\begin{array}{l} -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right)_1}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right)_2} \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1^{+1}}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \\ -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \frac{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1^{+1}}{\left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \eta_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \end{array} \right] &
 \end{aligned}$$

The next message has the form:

$$\begin{aligned}
 & \mathbf{m}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \\
 & \leftarrow \exp \left\{ \int_{\mathbb{R}^3} \int_{\mathbb{R}_{\geq 0}} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\alpha} \rightarrow p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\boldsymbol{\alpha}) \times \mathbf{m}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) \right. \\
 & \quad \times \mathbf{m}_{\sigma_o^2 \rightarrow p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\sigma_o^2) \times \mathbf{m}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2) \\
 & \quad \left. \times \log p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) d\boldsymbol{\alpha} d\sigma_o^2 \right\} \\
 & \propto \exp \left[\int_{\mathbb{R}^3} \int_{\mathbb{R}_{\geq 0}} \left\{ \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \right) \right\} \right. \\
 & \quad \times \left. \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right) \right\} \right. \\
 & \quad \left. \times \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \begin{bmatrix} \frac{\alpha_1}{\sigma_o^2} (\mathbf{o}_{x_{\text{unobs},i}} + \alpha_0) \\ -\frac{\alpha_1^2}{2\sigma_o^2} \end{bmatrix} d\boldsymbol{\alpha} d\sigma_o^2 \right] \\
 & = \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \right. \\
 & \quad \times \left. \left[\begin{aligned} & \int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_o^2} p_{\text{IG}_{\text{nat}}}(\sigma_o^2; \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}) d\sigma_o^2 \\ & \times \left\{ \mathbf{o}_{x_{\text{unobs},i}} \int_{\mathbb{R}} \alpha_1 p_{\text{N}_{\text{nat}}}(\alpha_1; \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} \right. \right. \\ & \quad \left. \left. + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1}) d\alpha_1 + \int_{\mathbb{R}} \int_{\mathbb{R}} \alpha_0 \alpha_1 p_{\text{N}_{\text{nat}}}(\alpha_0, \alpha_1; \right. \right. \\ & \quad \left. \left. \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_0} \boldsymbol{\eta}_{p(\alpha_0) \rightarrow \alpha_0}, \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} \right. \right. \\ & \quad \left. \left. + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1}) d\alpha_0 d\alpha_1 \right\} \right. \\ & \quad \left. - \frac{1}{2} \int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_o^2} p_{\text{IG}_{\text{nat}}}(\sigma_o^2; \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}) d\sigma_o^2 \right. \\ & \quad \left. \times \int_{\mathbb{R}} \alpha_1^2 p_{\text{N}_{\text{nat}}}(\alpha_1; \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1}) d\alpha_1 \right] \Bigg\}.
 \end{aligned}$$

Using Primitives 7.4.6, 7.4, and 7.4.1, this gives

$$\mathbf{m}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \propto \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \left[\begin{array}{l} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \left\{ -\frac{1}{2} \mathbf{o}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right)_2} \right. \\ \left. + \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} ; 1, 2 \right) \right\} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{2 \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \mathcal{K} \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right) \end{array} \right\}.$$

Therefore, the corresponding natural parameter update is

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \leftarrow \left[\begin{array}{l} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \left\{ -\frac{1}{2} \mathbf{o}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right)_2} \right. \\ \left. + \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} ; 1, 2 \right) \right\} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{2 \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \mathcal{K} \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \alpha_1} + \boldsymbol{\eta}_{p(\alpha_1) \rightarrow \alpha_1} \right) \end{array} \right].$$

Moving on, we next look at

$$\begin{aligned} \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \mathbf{a}}(\mathbf{a}) &\leftarrow \exp \left\{ \sum_{k=1}^K \sum_{\omega_k=0}^1 \frac{1}{Z} \mathbf{m}_{\boldsymbol{\omega} \rightarrow p(\mathbf{a}|\boldsymbol{\omega})}(\boldsymbol{\omega}) \times \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) \right. \\ &\quad \left. \times \log p(\mathbf{a}|\boldsymbol{\omega}) d\boldsymbol{\omega} \right\} \\ &\propto \exp \left\{ \sum_{k=1}^K \sum_{\omega_k=0}^1 \sum_{i=1}^n \log(\omega_k) \left(\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} + \boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} \right) \right. \\ &\quad \left. \times a_{ik} \log(\omega_k) d\omega_k \right\}. \end{aligned}$$

Using Result 7.4.2 gives

$$\begin{aligned} \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \mathbf{a}}(\mathbf{a}) &\propto \exp \left\{ \sum_{k=1}^K \sum_{i=1}^n a_{ik} \left(\psi \left\{ \left(\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_k + 1 \right\} \right. \right. \\ &\quad \left. \left. - \psi \left[\sum_{l=1}^K \left\{ \left(\boldsymbol{\eta}_{p(\omega_k) \rightarrow \omega_k} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_l + 1 \right\} \right] \right) \right\}. \end{aligned}$$

Thus, the corresponding natural parameter is

$$\eta_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \leftarrow \left(\psi \left\{ \left(\eta_{p(\omega_k) \rightarrow \omega_k} + \eta_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_k + 1 \right\} - \psi \left[\sum_{k=1}^K \left\{ \left(\eta_{p(\omega_k) \rightarrow \omega_k} + \eta_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \right)_l + 1 \right\} \right] \right).$$

Next, we consider the message

$$\begin{aligned} \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) &\leftarrow \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \frac{1}{Z} \mathbf{m}_{\mathbf{a} \rightarrow p(\mathbf{a}|\boldsymbol{\omega})}(\mathbf{a}) \times \mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \mathbf{a}}(\mathbf{a}) \right. \\ &\quad \left. \times \log p(\mathbf{a}|\boldsymbol{\omega}) \right\} \\ &\propto \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 a_{ik} \left(\eta_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}} + \eta_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \right) \right. \\ &\quad \left. \times \log(\omega_k) a_{ik} \right\}. \end{aligned}$$

The use of Result 7.4.1 gives

$$\mathbf{m}_{p(\mathbf{a}|\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) \propto \exp \left[\sum_{i=1}^n \sum_{k=1}^K \log(\omega_k) \left\{ \exp \left(\eta_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}} + \eta_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \right) \right\} \right].$$

This means that the update for the natural parameter vector is

$$\eta_{p(a_{ik}|\omega_k) \rightarrow \omega_k} \leftarrow \exp \left(\eta_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow a_{ik}} + \eta_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} \right).$$

Next we derive the update for the factor to node message $\mathbf{m}_{p(\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega})$:

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\omega}) \rightarrow \boldsymbol{\omega}}(\boldsymbol{\omega}) &\leftarrow \exp \left\{ \frac{1}{Z} \log p(\boldsymbol{\omega}) \right\} \\ &\propto \exp \left\{ \sum_{k=1}^K \log(\omega_k) (\alpha_k - 1) \right\}. \end{aligned}$$

Therefore the update for the corresponding natural parameter vector is

$$\eta_{p(\omega_k) \rightarrow \omega_k} \leftarrow \alpha_k - 1.$$

Similarly, we next consider another constant update, whose corresponding message takes the form

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) &\leftarrow \exp \left\{ \frac{1}{Z} \log p(\boldsymbol{\beta}) \right\} \\ &\propto \exp \left\{ \left[\begin{array}{c} \boldsymbol{\beta} \\ \text{vec}(\boldsymbol{\beta}\boldsymbol{\beta}^\top) \end{array} \right]^\top \left[\begin{array}{c} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right) \end{array} \right] \right\}, \end{aligned}$$

and so the natural parameter vector takes on the constant value

$$\boldsymbol{\eta}_{p(\beta) \rightarrow \beta} \leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_\beta^2} \mathbf{I}_3 \right) \end{bmatrix}.$$

Following on, we next consider

$$\begin{aligned} \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) &\leftarrow \exp \left\{ \int_{\mathbb{R}^{n_{\text{unobs}}}} \int_{\mathbb{R}_{\geq 0}} \frac{1}{Z} \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\mathbf{x}_{\text{unobs}}) \right. \\ &\quad \times \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \\ &\quad \times \mathbf{m}_{\sigma_o^2 \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\sigma_o^2) \\ &\quad \times \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2) \\ &\quad \left. \times \log p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) d\mathbf{x}_{\text{unobs}} d\sigma_o^2 \right\} \\ &\propto \exp \left[\int_{\mathbb{R}^{n_{\text{unobs}}}} \int_{\mathbb{R}_{\geq 0}} \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \boldsymbol{\alpha}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow x_{\text{unobs},i} \right. \right. \\ &\quad \left. \left. + \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \right) \right\} \\ &\quad \times \left[\begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow \sigma_o^2} \right) \right] \\ &\quad \times \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \begin{bmatrix} \frac{1}{\sigma_o^2} \tilde{\mathbf{X}}^\top \mathbf{o} \\ -\frac{1}{2\sigma_o^2} \text{vec}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \end{bmatrix} d\mathbf{x}_{\text{unobs}} d\sigma_o^2 \right] \\ &= \exp \left\{ \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \right. \\ &\quad \times \left. \begin{bmatrix} \int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_o^2} p_{\text{IG}_{\text{nat}}}(\sigma_o^2; \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow \sigma_o^2}) \\ \quad \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\tilde{\mathbf{X}})^\top \mathbf{o} \\ -\frac{1}{2} \int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_o^2} p_{\text{IG}_{\text{nat}}}(\sigma_o^2; \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2 | a_o) \rightarrow \sigma_o^2}) \\ \quad \times \text{vec} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right\} \end{bmatrix} \right\} \end{aligned}$$

$$= \exp \left\{ \left[\begin{array}{c} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{array} \right]^\top \left[\begin{array}{c} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}})^\top \mathbf{o} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{2 \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \text{vec} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right\} \end{array} \right] \right\}$$

where

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}) = \left[\begin{array}{cc} \mathbf{1}_{n_{\text{obs}}} & \mathbf{x}_{\text{obs}} \\ \mathbf{1}_{n_{\text{unobs}}} & -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \end{array} \right],$$

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) =$$

$$\left[\begin{array}{c|c} n & \mathbf{1}_{n_{\text{obs}}}^\top \mathbf{x}_{\text{obs}} \\ \hline -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} & \begin{array}{c} -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \\ \|\mathbf{x}_{\text{obs}}\|^2 + \\ \left\| -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \right\|^2 \\ -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \end{array} \end{array} \right]$$

and

$$\boldsymbol{\eta}^{\boxtimes} \equiv [\boldsymbol{\eta}_1^{\boxtimes} \boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes}]^\top \equiv \boldsymbol{\eta}_p\{\mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}} + \boldsymbol{\eta}_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}} + \boldsymbol{\eta}_p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}.$$

Therefore the corresponding natural parameter update is

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} \leftarrow \left[\begin{array}{c} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}})^\top \mathbf{o} \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1^{+1}}{2 \left(\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \\ \times \text{vec} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right\} \end{array} \right].$$

Next we have

$$\begin{aligned} \mathbf{m}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) &\leftarrow \exp \left\{ \frac{1}{2} \log p(\boldsymbol{\alpha}) \right\} \\ &\propto \exp \left\{ \left[\begin{array}{c} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{array} \right]^\top \left[\begin{array}{c} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_\alpha^2} \mathbf{I}_2 \right) \end{array} \right] \right\} \end{aligned}$$

Therefore the natural parameter update stays constant at:

$$\boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \leftarrow \begin{bmatrix} \mathbf{0} \\ -\frac{1}{2} \text{vec} \left(\frac{1}{\sigma_o^2} \mathbf{I}_2 \right) \end{bmatrix}.$$

The next message we look at is from $p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)$ to σ_o^2 :

$$\begin{aligned} & \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2) \\ & \leftarrow \exp \left\{ \int_{\mathbb{R}^2} \int_{\mathbb{R}^{n_{\text{unobs}}}} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\alpha} \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\boldsymbol{\alpha}) \times \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) \right. \\ & \quad \times \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}(\mathbf{x}_{\text{unobs}}) \times \mathbf{m}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \\ & \quad \left. \times \log p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) d\boldsymbol{\alpha}, d\mathbf{x}_{\text{unobs}} \right\} \\ & \propto \exp \left[\int_{\mathbb{R}^2} \int_{\mathbb{R}^{n_{\text{unobs}}}} \left\{ \begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha}\boldsymbol{\alpha}^\top) \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}} \right) \right\} \right. \\ & \quad \times \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{x} | \boldsymbol{\alpha}, \boldsymbol{\mu}^x, (\sigma^x)^2) \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} \right. \right. \\ & \quad \left. \left. + \boldsymbol{\eta}_{p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \right) \right\} \times \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \begin{bmatrix} -n/2 \\ -\frac{1}{2} \|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2 \end{bmatrix} d\boldsymbol{\alpha} d\mathbf{x}_{\text{unobs}} \Bigg] \\ & = \exp \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \begin{bmatrix} -n/2 \\ -\frac{1}{2} E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2) \end{bmatrix} \right\} \end{aligned}$$

where

$$\begin{aligned} & E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2) = \\ & \left\| \mathbf{o}_{x_{\text{obs}}} + \frac{1}{2} \tilde{\mathbf{X}}_{x_{\text{obs}}} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\star) \right\}^{-1} \boldsymbol{\eta}_1^\star \right\|^2 - \frac{1}{2} \text{tr} \left[\tilde{\mathbf{X}}_{x_{\text{obs}}}^\top \tilde{\mathbf{X}}_{x_{\text{obs}}} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\star) \right\}^{-1} \right] \\ & + \|\mathbf{o}_{x_{\text{unobs}}}\|^2 + \mathbf{o}_{x_{\text{unobs}}}^\top E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}_{x_{\text{unobs}}}) \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\star) \right\}^{-1} \boldsymbol{\eta}_1^\star \\ & + \frac{1}{4} \text{tr} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}_{x_{\text{unobs}}}^\top \tilde{\mathbf{X}}_{x_{\text{unobs}}}) \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\star) \right\}^{-1} \left[\boldsymbol{\eta}_1^\star (\boldsymbol{\eta}_1^\star)^\top \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\star) \right\}^{-1} \right. \right. \\ & \quad \left. \left. - 2\mathbf{I} \right] \right\}, \end{aligned}$$

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}_{x_{\text{unobs}}}) = \left[\mathbf{1}_{n_{\text{unobs}}} \quad -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\boxtimes) \right\}^{-1} \boldsymbol{\eta}_1^\boxtimes \right],$$

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\tilde{\mathbf{X}}_{x_{\text{unobs}}}^\top \tilde{\mathbf{X}}_{x_{\text{unobs}}}) =$$

$$\left[\begin{array}{c|c} n_{\text{unobs}} & -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\boxtimes) \right\}^{-1} \boldsymbol{\eta}_1^\boxtimes \\ \hline -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\boxtimes) \right\}^{-1} \boldsymbol{\eta}_1^\boxtimes & \left\| -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\boxtimes) \right\}^{-1} \boldsymbol{\eta}_1^\boxtimes \right\|^2 \\ & -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^\boxtimes) \right\}^{-1} \end{array} \right],$$

and

$$\boldsymbol{\eta}^\star \equiv [\boldsymbol{\eta}_1^\star \ \boldsymbol{\eta}_{2,\text{vec}}^\star]^\top \equiv \boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha}) \rightarrow \boldsymbol{\alpha}}.$$

Therefore, the update for the corresponding natural parameter is

$$\boldsymbol{\eta}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2) \leftarrow \begin{bmatrix} -n/2 \\ -\frac{1}{2} E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\|\mathbf{o} - \tilde{\mathbf{X}}\boldsymbol{\alpha}\|^2) \end{bmatrix}.$$

Following on, the next message is

$$\begin{aligned} \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}(\sigma_o^2) &\leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}} \frac{1}{Z} \mathbf{m}_{a_o \rightarrow p(\sigma_o^2|a_o)}(a_o) \times \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o}(a_o) \right. \\ &\quad \left. \times \log p(\sigma_o^2|a_o) da_o \right\} \\ &\propto \exp \left\{ \int_{\mathbb{R}_{\geq 0}} \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top (\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o}) \right. \\ &\quad \left. \times \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -1/a_o \end{bmatrix} da_o \right\} \\ &= \exp \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -\int_{\mathbb{R}_{\geq 0}} \frac{1}{a_o} p_{\text{IG}_{\text{nat}}}(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o}) da_o \end{bmatrix} \right\}. \end{aligned}$$

With use of Primitive 7.4.6, this gives

$$\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}(\sigma_o^2) \propto \exp \left\{ \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -\frac{(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o})_1 + 1}{(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o})_2} \end{bmatrix} \right\},$$

therefore the corresponding natural parameter update is

$$\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \leftarrow \begin{bmatrix} -3/2 \\ -\frac{(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o})_1 + 1}{(\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} + \boldsymbol{\eta}_{p(a_o) \rightarrow a_o})_2} \end{bmatrix}.$$

Next, we look at

$$\begin{aligned}
 \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o}(a_o) &\leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}} \frac{1}{Z} \mathbf{m}_{\sigma_o^2 \rightarrow p(\sigma_o^2|a_o)}(\sigma_o^2) \times \mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}(\sigma_o^2) \right. \\
 &\quad \left. \times \log p(\sigma_o^2|a_o) d\sigma_o^2 \right\} \\
 &\propto \exp \left\{ \int_{\mathbb{R}_{\geq 0}} \begin{bmatrix} \log(\sigma_o^2) \\ 1/\sigma_o^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right) \right. \\
 &\quad \left. \times \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2} \\ -1/\sigma_o^2 \end{bmatrix} d\sigma_o^2 \right\} \\
 &= \exp \left\{ \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2} \\ -\int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_o^2} p_{\text{IG}_{\text{nat}}} \left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right) d\sigma_o^2 \end{bmatrix} \right\}.
 \end{aligned}$$

Once again, making use of Primitive 7.4.6, we get

$$\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o}(a_o) \propto \exp \left\{ \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -\frac{\left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \end{bmatrix} \right\},$$

and so, the corresponding natural parameter update is

$$\boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow a_o} \leftarrow \begin{bmatrix} -3/2 \\ -\frac{\left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\sigma|a_o, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2} + \boldsymbol{\eta}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2} \right)_2} \end{bmatrix}.$$

Next, we have

$$\begin{aligned}
 \mathbf{m}_{p(a_o) \rightarrow a_o}(a_o) &\leftarrow \exp \left\{ \frac{1}{Z} \log(a_o) \right\} \\
 &\propto \exp \left\{ \begin{bmatrix} \log(a_o) \\ 1/a_o \end{bmatrix}^\top \begin{bmatrix} -3/2 \\ -1/A_o^2 \end{bmatrix} \right\},
 \end{aligned}$$

thus the corresponding natural parameter vector has the constant value

$$\boldsymbol{\eta}_{p(a_o) \rightarrow a_o} \leftarrow \begin{bmatrix} -3/2 \\ -1/A_o^2 \end{bmatrix}.$$

Next, we have

$$\begin{aligned}
 & \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\boldsymbol{\sigma}^x)^2 \{(\boldsymbol{\sigma}^x)^2\}} \\
 & \leftarrow \exp \left[\sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \int_{\mathbb{R}^{n_{\text{unobs}}}} \int_{\mathbb{R}^K} \frac{1}{Z} \mathbf{m}_{\mathbf{a} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}}(\mathbf{a}) \right. \\
 & \quad \times \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mathbf{a}}(\mathbf{a}) \\
 & \quad \times \mathbf{m}_{\mathbf{x}_{\text{unobs}} \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}}(\mathbf{x}_{\text{unobs}}) \times \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \\
 & \quad \times \mathbf{m}_{\boldsymbol{\mu}^x \rightarrow p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\}}(\boldsymbol{\mu}^x) \times \mathbf{m}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \boldsymbol{\mu}^x}(\boldsymbol{\mu}^x) \\
 & \quad \left. \times \log p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} d\mathbf{a} d\mathbf{x}_{\text{unobs}} d\boldsymbol{\mu}^x \right] \\
 & \propto \exp \left[\sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \int_{\mathbb{R}^{n_{\text{unobs}}}} \int_{\mathbb{R}^K} \left\{ a_{ik} \left(\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}} \right) \right\} \right. \\
 & \quad \times \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow x_{\text{unobs},i}} \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \right) \right\} \times \left\{ \sum_{k=1}^K \begin{bmatrix} \mu_k^x \\ (\mu_k^x)^2 \end{bmatrix}^\top \left(\boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x} \right) \right\} \left\{ \sum_{i=1}^n \sum_{k=1}^K \begin{bmatrix} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix}^\top \begin{bmatrix} -\frac{1}{2}a_{ik} \\ -\frac{1}{2}a_{ik}(x_i - \mu_k^x)^2 \end{bmatrix} \right\} \\
 & \quad \left. d\mathbf{a} d\mathbf{x}_{\text{unobs}} d\boldsymbol{\mu}^x \right] \\
 & \propto \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \begin{bmatrix} \log \{(\sigma_k^x)^2\} \\ 1/(\sigma_k^x)^2 \end{bmatrix}^\top \right. \\
 & \quad \left. \left[\begin{array}{l} -\frac{1}{2} \sum_{a_{ik}=0}^1 a_{ik} \times a_{ik} \times \left(\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}} \right) \\ -\frac{1}{2} \times a_{ik} \times \left(\boldsymbol{\eta}_{p(a_{ik}|\omega_k) \rightarrow a_{ik}} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik}} \right) \\ \times \int_{\mathbb{R}} \int_{\mathbb{R}} (x_i - \mu_k)^2 p_{\text{Nnat}}(x_i, \mu_k; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} \\ + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow x_{\text{unobs},i}} + \boldsymbol{\eta}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}} \\ \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x + \boldsymbol{\eta}_{p(\mu_k^x) \rightarrow \mu_k^x} \end{array} d\mathbf{x}_{\text{unobs},i} d\mu_k^x \right] \right\}.
 \end{aligned}$$

This implies that the (a_{i1}, \dots, a_{iK}) follow the

$$\begin{aligned}
 & \mathbf{M}_{\text{nat}} \left\{ 1; \left(\boldsymbol{\eta}_{p(a_i|\boldsymbol{\omega}) \rightarrow \mathbf{a}_i} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mathbf{a}_i} \right)_1, \dots, \right. \\
 & \quad \left. \left(\boldsymbol{\eta}_{p(a_i|\boldsymbol{\omega}) \rightarrow \mathbf{a}_i} + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mathbf{a}_i} \right)_K \right\}
 \end{aligned}$$

distribution and making use of Result 7.4.1 and Primitive 7.4.2, we get

$$\begin{aligned} & \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}^x)^2 \{ (\boldsymbol{\sigma}^x)^2 \} \propto \\ & \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \left[\log \left\{ (\sigma_k^x)^2 \right\} \right]^\top \right. \\ & \quad \times \left. \left[\begin{array}{l} -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \mathcal{I} \left(\eta_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} \right. \\ \quad \left. + \eta_p(\sigma | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}, \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x \right. \\ \quad \left. \left. + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right) \right] \right\}. \end{array} \right. \end{aligned}$$

Therefore, the update for the corresponding natural parameter vector is

$$\begin{aligned} & \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}^x)^2 \leftarrow \\ & \left[\begin{array}{l} -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ -\frac{1}{2} \exp \left(\eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \right) \\ \times \mathcal{I} \left(\eta_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} \right. \\ \quad \left. + \eta_p(\sigma | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}, \eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x \right. \\ \quad \left. \left. + \eta_p(\mu_k^x) \rightarrow \mu_k^x \right) \right] \end{array} \right. \end{aligned}$$

Next we look at

$$\begin{aligned} & \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{a}(\mathbf{a}) \\ & \leftarrow \exp \left[\int_{\mathbb{R}^K} \int_{\mathbb{R}^K} \int_{\mathbb{R}^{n_{\text{unobs}}}} \frac{1}{Z} \mathbf{m} \left((\boldsymbol{\sigma}^x)^2 \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \{ (\boldsymbol{\sigma}^x)^2 \} \right. \right. \\ & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}^x)^2 \{ (\boldsymbol{\sigma}^x)^2 \} \times \mathbf{m}_{\boldsymbol{\mu}^x} \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} (\boldsymbol{\mu}^x) \\ & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \boldsymbol{\mu}^x (\boldsymbol{\mu}^x) \times \mathbf{m}_{\mathbf{x}_{\text{unobs}}} \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} (\mathbf{x}_{\text{unobs}}) \\ & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{x}_{\text{unobs}} (\mathbf{x}_{\text{unobs}}) \times \log p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \\ & \quad \left. \left. d(\boldsymbol{\sigma}^x)^2 d\boldsymbol{\mu}^x d\mathbf{x}_{\text{unobs}} \right] \end{aligned}$$

$$\begin{aligned}
 & \propto \exp \left[\int_{\mathbb{R}^K} \int_{\mathbb{R}^K} \int_{\mathbb{R}^{n_{\text{unobs}}}} \left\{ \sum_{k=1}^K \left[\log \left\{ \frac{(\sigma_k^x)^2}{1/(\sigma_k^x)^2} \right\} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 \right. \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right\} \times \left\{ \sum_{k=1}^K \left[\frac{\mu_k^x}{(\mu_k^x)^2} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p (\mu_k^x \rightarrow \mu_k^x) \right) \right\} \times \left\{ \sum_{i=1}^{n_{\text{unobs}}} \sum_{k=1}^K \left[\frac{x_{\text{unobs},i}}{x_{\text{unobs},i}^2} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p (\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p (\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i} \right) \right\} \\
 & \quad \times \sum_{i=1}^n \sum_{k=1}^K a_{ik} \left(\log \left[\{ 2\pi (\sigma_k^x)^2 \}^{-1/2} \exp \left\{ -\frac{1}{2} (x_i - \mu_k^x)^2 / (\sigma_k^x)^2 \right\} \right] \right) \\
 & \quad d(\boldsymbol{\sigma}^x)^2 d\boldsymbol{\mu}^x d\mathbf{x}_{\text{unobs}} \\
 & \propto \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K a_{ik} \left(-\frac{1}{2} \left[\int_{\mathbb{R}_{\geq 0}} \log \{ (\sigma_k^x)^2 \} p_{\text{IG}_{\text{nat}}} \left((\sigma_k^x)^2 ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 \right. \right. \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) d(\sigma_k^x)^2 - \int_{\mathbb{R}_{\geq 0}} \frac{1}{(\sigma_k^x)^2} p_{\text{IG}_{\text{nat}}} \left((\sigma_k^x)^2 ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) d(\sigma_k^x)^2 \times \int_{\mathbb{R}} \int_{\mathbb{R}} (x_i - \mu_k^x)^2 p_{\text{N}_{\text{nat}}} (x_i, \mu_k^x ; \right. \\
 & \quad \left. \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p (\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p (\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}, \right. \\
 & \quad \left. \left. \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \boldsymbol{\eta}_p (\mu_k^x \rightarrow \mu_k^x) \right) d x_{\text{unobs},i} d \mu_k^x \right] \right\}.
 \end{aligned}$$

Making use of Primitives 7.4.2, 7.4.6 and 7.4.7, this gives

$$\begin{aligned}
 & \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{a} (\mathbf{a}) \propto \\
 & \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K a_{ik} \left(-\frac{1}{2} \left[\log \left\{ - \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right\} \right. \right. \\
 & \quad \left. \left. - \psi \left\{ - \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) \right\} - 1 \right\} \right. \\
 & \quad \left. - \frac{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1 + 1}{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \right\} \\
 & \quad \times \mathcal{I} \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} \boldsymbol{\eta}_p (\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p (\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i}, \right. \\
 & \quad \left. \left. \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \boldsymbol{\eta}_p (\mu_k^x \rightarrow \mu_k^x) \right) \right] \right\}.
 \end{aligned}$$

Recognising that \mathbf{a} follows a series of n Multinomial distributions, the update for the

corresponding natural parameter vector is

$$\begin{aligned}
 & \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} \leftarrow \\
 & -\frac{1}{2} \left[\log \left\{ - \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 + \boldsymbol{\eta}_p \{ (\boldsymbol{\sigma}_k^x)^2 | a_k^x \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right)_2 \right\} \right. \\
 & \left. - \psi \left\{ - \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 + \boldsymbol{\eta}_p \{ (\boldsymbol{\sigma}_k^x)^2 | a_k^x \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right)_1 - 1 \right\} \right. \\
 & \left. - \frac{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 + \boldsymbol{\eta}_p \{ (\boldsymbol{\sigma}_k^x)^2 | a_k^x \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right)_1 + 1}{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 + \boldsymbol{\eta}_p \{ (\boldsymbol{\sigma}_k^x)^2 | a_k^x \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right)_2} \right] \\
 & \times \mathcal{I} \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} \boldsymbol{\eta}_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \boldsymbol{\sigma}_\varepsilon^2) \rightarrow x_{\text{unobs},i} \right. \\
 & \left. + \boldsymbol{\eta}_p(o | \boldsymbol{\alpha}, \mathbf{x}, \boldsymbol{\sigma}_o^2) \rightarrow x_{\text{unobs},i}, \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mu_k^x + \boldsymbol{\eta}_p(\mu_k^x) \rightarrow \mu_k^x \right).
 \end{aligned}$$

Next,

$$\begin{aligned}
 & \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \boldsymbol{\mu}^x (\boldsymbol{\mu}^x) \\
 & \leftarrow \exp \left\{ \int_{\mathbb{R}_{\geq 0}} \sum_{i=1}^n \sum_{k=1}^K \sum_{\mathbf{a}=0}^1 \int_{\mathbb{R}^{n_{\text{unobs}}}} \frac{1}{Z} \mathbf{m} (\boldsymbol{\sigma}^x)^2 \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \{ (\boldsymbol{\sigma}^x)^2 \} \right. \\
 & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}^x)^2 \{ (\boldsymbol{\sigma}^x)^2 \} \times \mathbf{m}_p \mathbf{a} \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} (\mathbf{a}) \\
 & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{a} (\mathbf{a}) \times \mathbf{m}_p \mathbf{x}_{\text{unobs}} \rightarrow p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} (\mathbf{x}_{\text{unobs}}) \\
 & \quad \times \mathbf{m}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow \mathbf{x}_{\text{unobs}} (\mathbf{x}_{\text{unobs}}) \times \log p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \\
 & \quad \left. d (\boldsymbol{\sigma}^x)^2 d \mathbf{a} d \mathbf{x}_{\text{unobs}} \right\} \\
 & \propto \left[\int_{\mathbb{R}_{\geq 0}} \sum_{i=1}^n \sum_{k=1}^K \sum_{a_{ik}=0}^1 \int_{\mathbb{R}^{n_{\text{unobs}}}} \left\{ \sum_{k=1}^K \left[\log \{ (\boldsymbol{\sigma}_k^x)^2 \} \right] \right\}^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p \{ (\boldsymbol{\sigma}_k^x)^2 | a_k^x \} \rightarrow (\boldsymbol{\sigma}_k^x)^2 \right) \times \left\{ \sum_{i=1}^n \sum_{k=1}^K a_{ik} \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow a_{ik} + \boldsymbol{\eta}_p(a_{ik} | \omega_k) \rightarrow a_{ik} \right) \right\} \right. \\
 & \quad \times \left\{ \sum_{i=1}^{n_{\text{unobs}}} \left[\begin{array}{c} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{array} \right]^\top \left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2 \} \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \boldsymbol{\sigma}_\varepsilon^2) \rightarrow x_{\text{unobs},i} \right. \right. \\
 & \quad \left. \left. + \boldsymbol{\eta}_p(o | \boldsymbol{\alpha}, \mathbf{x}, \boldsymbol{\sigma}_o^2) \rightarrow x_{\text{unobs},i} \right) \right\} \times \left\{ \sum_{i=1}^n \sum_{k=1}^K \left[\begin{array}{c} \mu_k^x \\ (\mu_k^x)^2 \end{array} \right]^\top \left[\begin{array}{c} a_{ik} (x_i / (\boldsymbol{\sigma}_k^x)^2) \\ -\frac{1}{2} a_{ik} / (\boldsymbol{\sigma}_k^x)^2 \end{array} \right] \right\} \\
 & \quad \left. d (\boldsymbol{\sigma}^x)^2 d \mathbf{a} d \mathbf{x}_{\text{unobs}} \right]
 \end{aligned}$$

$$\propto \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \left[\begin{array}{l} \mu_k^x \\ (\mu_k^x)^2 \end{array} \right]^\top \right. \\ \left. \times \left[\begin{array}{l} \sum_{a_{ik}=0}^1 a_{ik} \times a_{ik} \times \left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow a_{ik} + \eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} \right) \\ \times \int_{\mathbb{R}} x_i p_{\text{Nnat}} \left(x_i ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow x_{\text{unobs},i} \right. \\ \left. + \boldsymbol{\eta}_p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_p(\mathbf{o} | \boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow x_{\text{unobs},i} \right) d x_{\text{unobs},i} \\ \times \int_{\mathbb{R}_{\geq 0}} \frac{1}{(\sigma_k^x)^2} p_{\text{IGnat}} \left((\sigma_k^x)^2 ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 \right. \\ \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) d (\sigma_k^x)^2 \\ \\ - \frac{1}{2} \sum_{a_{ik}=0}^1 a_{ik} \times a_{ik} \times \left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow a_{ik} + \eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} \right) \\ \times \int_{\mathbb{R}_{\geq 0}} \frac{1}{(\sigma_k^x)^2} p_{\text{IGnat}} \left((\sigma_k^x)^2 ; \boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 \right. \\ \left. + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right) d (\sigma_k^x)^2 \end{array} \right] \right\}.$$

Using Primitive 7.4.6 and Result 7.4.1, we get

$$\propto \exp \left\{ \sum_{i=1}^n \sum_{k=1}^K \left[\begin{array}{l} \mu_k^x \\ (\mu_k^x)^2 \end{array} \right]^\top \right. \\ \left. \times \left[\begin{array}{l} -\frac{1}{2} \exp \left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow a_{ik} + \eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} \right) \\ \times \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{x}) \right\}_i \\ \times \frac{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1 + 1}{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \\ \\ - \frac{1}{2} \exp \left(\eta_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow a_{ik} + \eta_p(a_{ik} | \omega_k) \rightarrow a_{ik} \right) \\ \times \frac{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_1 + 1}{\left(\boldsymbol{\eta}_p \{ \mathbf{x} | \mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2 \} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p \{ (\sigma_k^x)^2 | a_k^x \} \rightarrow (\sigma_k^x)^2 \right)_2} \end{array} \right] \right\}$$

where

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{x}) \equiv \begin{bmatrix} \mathbf{x}_{\text{obs}} \\ -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes} \right) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \end{bmatrix}$$

therefore the corresponding natural parameter vector is

$$\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow \mu_k^x \leftarrow \begin{bmatrix} -\frac{1}{2}\exp\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik} + \boldsymbol{\eta}_p(a_{ik}|\omega_k) \rightarrow a_{ik}\right) \\ \times \left\{E_q^{nat}(\mathbf{x})\right\}_i \\ \times \frac{\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2\right)_1 + 1}{\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2\right)_2} \\ -\frac{1}{2}\exp\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow a_{ik} + \boldsymbol{\eta}_p(a_{ik}|\omega_k) \rightarrow a_{ik}\right) \\ \times \frac{\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2\right)_1 + 1}{\left(\boldsymbol{\eta}_p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\boldsymbol{\sigma}^x)^2\} \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_p\{(\sigma_k^x)^2|a_k^x\} \rightarrow (\sigma_k^x)^2\right)_2} \end{bmatrix}.$$

Continuing on, we next consider

$$\begin{aligned} & \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}}) \\ & \leftarrow \exp\left\{\int_{\mathbb{R}^3} \int_{\mathbb{R}_{\geq 0}} \frac{1}{Z} \mathbf{m}_{\boldsymbol{\beta} \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)}(\boldsymbol{\beta}) \times \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}}(\boldsymbol{\beta}) \right. \\ & \quad \times \mathbf{m}_{\sigma_\varepsilon^2 \rightarrow p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)}(\sigma_\varepsilon^2) \times \mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2) \\ & \quad \left. \times \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) d\boldsymbol{\beta} d\sigma_\varepsilon^2\right\} \\ & \propto \exp\left[\int_{\mathbb{R}^3} \int_{\mathbb{R}_{\geq 0}} \left\{\left[\begin{array}{c} \boldsymbol{\beta} \\ \text{vec}(\boldsymbol{\beta}\boldsymbol{\beta}^\top) \end{array}\right]^\top \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}}\right)\right\} \right. \\ & \quad \times \left\{\left[\begin{array}{c} \log(\sigma_\varepsilon^2) \\ 1/\sigma_\varepsilon^2 \end{array}\right]^\top \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}\right)\right\} \\ & \quad \left. \times \sum_{i=1}^{n_{\text{unobs}}} \left[\begin{array}{c} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{array}\right]^\top \left[\begin{array}{c} \frac{1}{\sigma_\varepsilon^2} (\beta_x \mathbf{y}_{x_{\text{unobs},i}} - \beta_0 \beta_x - \beta_c \beta_x \mathbf{c}_{x_{\text{unobs},i}}) \\ -\frac{\beta_x}{2\sigma_\varepsilon^2} \end{array}\right] d\boldsymbol{\beta} d\sigma_\varepsilon^2\right] \end{aligned}$$

$$\begin{aligned}
 &= \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \right. \\
 &\quad \left. \begin{aligned}
 &\int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_\varepsilon^2} p_{\text{IG}_{\text{nat}}} \left(\sigma_\varepsilon^2; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right) d\sigma_\varepsilon^2 \\
 &\times \left\{ \mathbf{y}_{x_{\text{unobs},i}} \int_{\mathbb{R}} \beta_x p_{\text{N}_{\text{nat}}} \left(\beta_x; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right) d\beta_x \right. \\
 &\quad \left. - \int_{\mathbb{R}^2} \beta_0 \beta_x p_{\text{N}_{\text{nat}}} \left(\beta_0, \beta_x; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_0 + \boldsymbol{\eta}_{p(\beta_0)} \rightarrow \beta_0, \right. \right. \\
 &\quad \left. \left. \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right) d\beta_0 d\beta_x \right. \\
 &\quad \left. - c_{x_{\text{unobs},i}} \int_{\mathbb{R}^2} \beta_c \beta_x p_{\text{N}_{\text{nat}}} \left(\beta_c, \beta_x; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_c + \boldsymbol{\eta}_{p(\beta_c)} \rightarrow \beta_c, \right. \right. \\
 &\quad \left. \left. \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right) d\beta_c d\beta_x \right\} \\
 &\quad \left. - \frac{1}{2} \int_{\mathbb{R}_{\geq 0}} \frac{1}{\sigma_\varepsilon^2} p_{\text{IG}_{\text{nat}}} \left(\sigma_\varepsilon^2; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right) d\sigma_\varepsilon^2 \right. \\
 &\quad \left. \times \int_{\mathbb{R}} \beta_x p_{\text{IG}_{\text{nat}}} \left(\beta_x; \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right) d\beta_x \right. \\
 &\quad \left. \right\}
 \end{aligned}
 \right.
 \end{aligned}$$

With the use of Primitives 7.4.5 and 7.4.6, this gives

$$\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \mathbf{x}_{\text{unobs}}}(\mathbf{x}_{\text{unobs}})$$

$$\begin{aligned}
 &\propto \exp \left\{ \sum_{i=1}^{n_{\text{unobs}}} \begin{bmatrix} x_{\text{unobs},i} \\ x_{\text{unobs},i}^2 \end{bmatrix}^\top \right. \\
 &\quad \left. \begin{aligned}
 &\frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right)_2} \\
 &\times \left\{ -\frac{1}{2} \mathbf{y}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right)_2} \right. \\
 &\quad \left. - \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta + \boldsymbol{\eta}_{p(\beta)} \rightarrow \beta; 1, 3 \right) \right. \\
 &\quad \left. - c_{x_{\text{unobs},i}} \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta + \boldsymbol{\eta}_{p(\beta)} \rightarrow \beta; 2, 3 \right) \right\} \\
 &\quad \frac{1}{4} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \sigma_\varepsilon^2 + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon)} \rightarrow \sigma_\varepsilon^2 \right)_2} \\
 &\quad \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)} \rightarrow \beta_x + \boldsymbol{\eta}_{p(\beta_x)} \rightarrow \beta_x \right)_2} \\
 &\quad \left. \right\}
 \end{aligned}
 \right.
 \end{aligned}$$

and so the corresponding natural parameter update is

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow x_{\text{unobs},i}} \leftarrow \left[\begin{aligned} & \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_2} \\ & \times \left\{ -\frac{1}{2} \mathbf{y}_{x_{\text{unobs},i}} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} \right)_2} \right. \\ & \quad - \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} ; 1, 3 \right) \\ & \quad \left. - \mathbf{c}_{x_{\text{unobs},i}} \mathcal{J} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} + \boldsymbol{\eta}_{p(\boldsymbol{\beta}) \rightarrow \boldsymbol{\beta}} ; 2, 3 \right) \right\} \\ & \frac{1}{4} \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_2} \\ & \times \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} \right)_1}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \beta_x} + \boldsymbol{\eta}_{p(\beta_x) \rightarrow \beta_x} \right)_2} \end{aligned} \right]$$

Next we consider the message from factor $p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2)$ to stochastic node $\boldsymbol{\beta}$. The derivation for this message is similar to that of $\mathbf{m}_{p(\boldsymbol{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha})$ and has the form:

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \boldsymbol{\beta}} \leftarrow \left[\begin{aligned} & \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_1^{+1}}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_2} \\ & \quad \times E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X})^\top \mathbf{y} \\ & \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_1^{+1}}{2 \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_2} \\ & \quad \times \text{vec} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X}^\top \mathbf{X}) \right\} \end{aligned} \right],$$

where

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X}) = \begin{bmatrix} \mathbf{1}_{n_{\text{obs}}} & \mathbf{c}_{\text{obs}} & \mathbf{x}_{\text{obs}} \\ \mathbf{1}_{n_{\text{unobs}}} & \mathbf{c}_{\text{unobs}} & -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\boldsymbol{\eta}_{2, \text{vec}}^{\boxtimes} \right) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \end{bmatrix}$$

and

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X}^\top \mathbf{X}) =$$

7.B. DERIVATION OF ALGORITHM 12

$$\left[\begin{array}{c|c|c} & & \mathbf{1}_{n_{\text{obs}}}^\top \mathbf{x}_{\text{obs}} \\ \hline n & \mathbf{1}_n^\top \mathbf{c} & -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} \\ & & \quad \times \eta_1^\boxtimes \\ \hline \mathbf{1}_n^\top \mathbf{c} & \|\mathbf{c}\|^2 & \mathbf{c}_{x_{\text{obs}}}^\top \mathbf{x}_{\text{obs}} \\ & & -\frac{1}{2} \mathbf{c}_{x_{\text{unobs}}}^\top \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} \\ & & \quad \times \eta_1^\boxtimes \\ \hline -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^\top \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} & \mathbf{c}_{x_{\text{obs}}}^\top \mathbf{x}_{\text{obs}} & \|\mathbf{x}_{\text{obs}}\|^2 + \\ & -\frac{1}{2} \mathbf{c}_{x_{\text{unobs}}}^\top \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} & -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} \eta_1^\boxtimes \|^2 \\ & \quad \times \eta_1^\boxtimes & -\frac{1}{2} \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\boxtimes \right) \right\}^{-1} \end{array} \right].$$

The last four messages we consider are $\mathbf{m}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2)$, $\mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2)$, $\mathbf{m}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon)$ and $\mathbf{m}_{p(a_\varepsilon) \rightarrow a_\varepsilon}(a_\varepsilon)$. The derivation of these messages is similar to that shown for $\mathbf{m}_{p(\mathbf{o}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2) \rightarrow \sigma_o^2}(\sigma_o^2)$, $\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow \sigma_o^2}(\sigma_o^2)$, $\mathbf{m}_{p(\sigma_o^2|a_o) \rightarrow a_o}(a_o)$ and $\mathbf{m}_{p(a_o) \rightarrow a_o}(a_o)$, respectively. Consequently, we only give their corresponding natural parameter updates:

$$\begin{aligned} \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2}(\sigma_\varepsilon^2) &\leftarrow \begin{bmatrix} -n/2 \\ -\frac{1}{2} E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2) \end{bmatrix}, \\ \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} &\leftarrow \begin{bmatrix} -3/2 \\ \frac{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} + \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} \right)_2} \end{bmatrix}, \\ \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow a_\varepsilon} &\leftarrow \begin{bmatrix} -3/2 \\ \frac{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_1 + 1}{\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_\varepsilon^2) \rightarrow \sigma_\varepsilon^2} + \boldsymbol{\eta}_{p(\sigma_\varepsilon^2|a_\varepsilon) \rightarrow \sigma_\varepsilon^2} \right)_2} \end{bmatrix}, \\ \boldsymbol{\eta}_{p(a_\varepsilon) \rightarrow a_\varepsilon} &\leftarrow \begin{bmatrix} -3/2 \\ -1/A_\varepsilon^2 \end{bmatrix}, \end{aligned}$$

where

$$\begin{aligned} E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2) &= \\ &\left\| \mathbf{y}_{x_{\text{obs}}} + \frac{1}{2} \mathbf{X}_{x_{\text{obs}}} \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\diamond \right) \right\}^{-1} \eta_1^\diamond \right\|^2 - \frac{1}{2} \text{tr} \left[\mathbf{X}_{x_{\text{obs}}}^\top \mathbf{X}_{x_{\text{obs}}} \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\diamond \right) \right\}^{-1} \right] \\ &+ \|\mathbf{y}_{x_{\text{unobs}}}\|^2 + \mathbf{y}_{x_{\text{unobs}}}^\top E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X}_{x_{\text{unobs}}}) \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\diamond \right) \right\}^{-1} \eta_1^\diamond \\ &+ \frac{1}{4} \text{tr} \left\{ E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}} (\mathbf{X}_{x_{\text{unobs}}}^\top \mathbf{X}_{x_{\text{unobs}}}) \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\diamond \right) \right\}^{-1} \left[\eta_1^\diamond (\eta_1^\diamond)^\top \left\{ \text{vec}^{-1} \left(\eta_{2,\text{vec}}^\diamond \right) \right\}^{-1} \right. \right. \\ &\left. \left. - 2\mathbf{I} \right] \right\}, \end{aligned}$$

7.B. DERIVATION OF ALGORITHM 12

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{X}_{x_{\text{unobs}}}) = \left[\mathbf{1}_{n_{\text{unobs}}} \quad \mathbf{c}_{x_{\text{unobs}}} \quad -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \right],$$

$$E_{q(\mathbf{x}_{\text{unobs}})}^{\text{nat}}(\mathbf{X}_{x_{\text{unobs}}}^{\top} \mathbf{X}_{x_{\text{unobs}}}) =$$

$$\begin{bmatrix} n_{\text{unobs}} & \mathbf{1}_{n_{\text{unobs}}}^{\top} \mathbf{c}_{x_{\text{unobs}}} & -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^{\top} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \times \boldsymbol{\eta}_1^{\boxtimes} \\ \mathbf{1}_{n_{\text{unobs}}}^{\top} \mathbf{c}_{x_{\text{unobs}}} & \|\mathbf{c}_{x_{\text{unobs}}}\|^2 & -\frac{1}{2} \mathbf{c}_{x_{\text{unobs}}}^{\top} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \times \boldsymbol{\eta}_1^{\boxtimes} \\ -\frac{1}{2} \mathbf{1}_{n_{\text{unobs}}}^{\top} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \times \boldsymbol{\eta}_1^{\boxtimes} & -\frac{1}{2} \mathbf{c}_{x_{\text{unobs}}}^{\top} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \times \boldsymbol{\eta}_1^{\boxtimes} & \left\| -\frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \boldsymbol{\eta}_1^{\boxtimes} \right\|^2 - \frac{1}{2} \left\{ \text{vec}^{-1}(\boldsymbol{\eta}_{2,\text{vec}}^{\boxtimes}) \right\}^{-1} \end{bmatrix},$$

$$\text{and } \boldsymbol{\eta}^{\diamond} \equiv [\boldsymbol{\eta}_1^{\diamond} \quad \boldsymbol{\eta}_{2,\text{vec}}^{\diamond}]^{\top} \equiv \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_{\varepsilon}^2)} \rightarrow \boldsymbol{\beta} + \boldsymbol{\eta}_{p(\boldsymbol{\beta})} \rightarrow \boldsymbol{\beta}.$$

7.B.4 Step 3 : Optimal q -densities

Once convergence is achieved from Steps 2(a) and (b) we get the optimal q -densities of each stochastic node through the following updates:

$$\boldsymbol{\eta}_{q(\boldsymbol{\alpha})}^* \leftarrow \boldsymbol{\eta}_{p(\boldsymbol{\alpha}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}^* \rightarrow \boldsymbol{\alpha} + \boldsymbol{\eta}_{p(\boldsymbol{\alpha})}^* \rightarrow \boldsymbol{\alpha}$$

$$\boldsymbol{\eta}_{q(\sigma_o^2)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_o^2|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}^* \rightarrow \sigma_o^2 + \boldsymbol{\eta}_{p(\sigma_o^2|a_o)}^* \rightarrow \sigma_o^2$$

$$\boldsymbol{\eta}_{q(a_o)}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_o^2|a_o)}^* \rightarrow a_o + \boldsymbol{\eta}_{p(a_o)}^* \rightarrow a_o.$$

For $1 \leq k \leq K$:

$$\boldsymbol{\eta}_{q(a_k^x)}^* \leftarrow \boldsymbol{\eta}_{p(a_k^x)}^* \rightarrow a_k^x + \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\}}^* \rightarrow a_k^x$$

$$\boldsymbol{\eta}_{q\{(\sigma_k^x)^2\}}^* \leftarrow \boldsymbol{\eta}_{p\{(\sigma_k^x)^2|a_k^x\}}^* \rightarrow (\sigma_k^x)^2 + \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}}^* \rightarrow (\sigma_k^x)^2$$

$$\boldsymbol{\eta}_{q(\mu_k^x)}^* \leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}}^* \rightarrow \mu_k^x + \boldsymbol{\eta}_{p(\mu_k^x)}^* \rightarrow \mu_k^x.$$

For $1 \leq i \leq n$ and $1 \leq k \leq K$:

$$\boldsymbol{\eta}_{q(a_{ik})}^* \leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}}^* \rightarrow a_{ik} + \boldsymbol{\eta}_{p(a_{ik}|\omega_k)}^* \rightarrow a_{ik}$$

$$\boldsymbol{\eta}_{q(\omega_k)}^* \leftarrow \boldsymbol{\eta}_{p(a_{ik}|\omega_k)}^* \rightarrow \omega_k + \boldsymbol{\eta}_{p(\omega_k)}^* \rightarrow \omega_k$$

$$\boldsymbol{\eta}_{q(\boldsymbol{\beta})}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_{\varepsilon}^2)}^* \rightarrow \boldsymbol{\beta} + \boldsymbol{\eta}_{p(\boldsymbol{\beta})}^* \rightarrow \boldsymbol{\beta}$$

$$\boldsymbol{\eta}_{q(\sigma_{\varepsilon}^2)}^* \leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_{\varepsilon}^2)}^* \rightarrow \sigma_{\varepsilon}^2 + \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon})}^* \rightarrow \sigma_{\varepsilon}^2$$

$$\boldsymbol{\eta}_{q(a_{\varepsilon})}^* \leftarrow \boldsymbol{\eta}_{p(\sigma_{\varepsilon}^2|a_{\varepsilon})}^* \rightarrow a_{\varepsilon} + \boldsymbol{\eta}_{p(a_{\varepsilon})}^* \rightarrow a_{\varepsilon}.$$

For $1 \leq i \leq n_{\text{unobs}}$:

$$\begin{aligned} \boldsymbol{\eta}_{q(x_{\text{unobs},i})}^* &\leftarrow \boldsymbol{\eta}_{p\{\mathbf{x}|\mathbf{a}, \boldsymbol{\mu}^x, (\sigma^x)^2\}}^* \rightarrow x_{\text{unobs},i} + \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x}, \sigma_{\varepsilon}^2)}^* \rightarrow x_{\text{unobs},i} \\ &+ \boldsymbol{\eta}_{p(\boldsymbol{\alpha}|\boldsymbol{\alpha}, \mathbf{x}, \sigma_o^2)}^* \rightarrow x_{\text{unobs},i}. \end{aligned}$$

Chapter 8

Conclusion

This thesis has aimed to increase the body of algorithms and results on MFVB inference in heteroscedastic and longitudinal regression. It has helped emphasise the importance of MFVB in both semiparametric and nonparametric regression settings and has demonstrated that MFVB provides a fast deterministic alternative to the stochastic MCMC with often little degradation in accuracy. Overall, we have developed extensions of previously used MFVB methodology for some model settings, whereas in others, we have developed original MFVB algorithms which had not previously been considered. Even though the simulations and real datasets used in this thesis have not exceeded over 10 megabytes in size, the idea behind the construction of these fast algorithms is to apply such methodology to larger and even more complex datasets that require further storage than what is available on a standard computer.

In Chapter 2, we extended the work of Al Kadiri *et al.* (2010) to obtain an MFVB algorithm catered to the marginal longitudinal semiparametric regression setting. This resulted in faster inference for the nutritional epidemiology study and it was clear that MFVB estimation of the model parameters in a simulation study was very high.

Chapter 3 considered the setting involving heteroscedastic semiparametric regression. This led to a modification of ordinary MFVB, called non-conjugate variational message passing which aids the incorporation of a heteroscedastic component without the need to stray from closed form algebraic expressions. Overall, the accuracy of the mean functions is very high and the accuracy of the variance functions, whilst not as accurate as the mean, are still producing considerably favorable accuracy scores.

In Chapter 4, we looked into the development of MFVB algorithms for three extensions

of the heteroscedastic setting considered in Chapter 3. We explore a real-time, bivariate predictor and additive model analogue of Chapter 3. Through the use of real-data examples, simulation studies and a real-time movie, we saw the time benefits of using such an algorithm.

Chapter 5 looks into the development of an MFVB algorithm for subject-specific curve models making use of longitudinal and multilevel structures. This chapter takes advantage of the structure of the model and looks into streamlining a naïve implementation of an MFVB algorithm, resulting in even faster inference.

In Chapter 6, we derived a fast MFVB algorithm to fitting a linear regression model with classical measurement error. The time savings here were evident, however extension of such methodology to arbitrarily large models involving measurement error would show time savings of higher magnitudes.

Chapter 7 discussed VMP, an alternative to MFVB, which produces the same approximation but with a different algebraic system. There, we derived the VMP updates for two models used throughout the thesis and provided the details of the methodology behind VMP by presenting a general VMP scheme.

It is generally not known ahead of time whether the MFVB approximation will produce an acceptable level of accuracy. The accuracy of MFVB varies depending on the type of application. As exhibited in this thesis, MFVB estimation of mean curves is very good but MFVB credible sets are generally underestimated. However, even with the shortcomings of poor estimation of credible sets, the speed and computational gains afforded by MFVB make it a good choice for complex models combined with large datasets.

Each chapter has shed light on the overall effectiveness of MFVB inference as a fast alternative to MCMC when time and/or storage is of concern. We have seen that the accuracy of MFVB varies depending on many details involved in the type of application. As indicated by the variety of models used, MFVB is of no doubt a dynamic area of statistical research where the potential to further contributions is profound.

Bibliography

- Aitkin, M. & Rocci, R. (2002). A general maximum likelihood analysis of measurement error in generalized linear models. *Statistics and Computing*, **12**, no. 2, 163–174.
- Al Kadiri, M., Carroll, R. & Wand, M. (2010). Marginal longitudinal semiparametric regression via penalized splines. *Statistics & probability letters*, **80**, no. 15, 1242–1252.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*, volume 1. springer New York.
- Bishop, C. M., Spiegelhalter, D. J. & Winn, J. (2003). *VIBES: A variational inference engine for Bayesian networks*. In advances in Neural Information Processing Systems, eds S. Becker, S. Thrun and K. Obermayer. Cambridge, MA: MIT Press.
- Braun, M. & McAuliffe, J. (2010). Variational inference for large-scale models of discrete choice. *Journal of the American Statistical Association*, **105**, no. 489, 324–335.
- Breiman, L. & Friedman, J. H. (1985). Estimating optimal transformations for multiple regression and correlation. *Journal of the American statistical Association*, **80**, no. 391, 580–598.
- Bugbee, B. D., Breidt, F. J. & van der Woerd, M. J. (2015). Laplace variational approximation for semiparametric regression in the presence of heteroskedastic errors. *Journal of Computational and Graphical Statistics*, in press.
- Carroll, R., Gail, M. & Lubin, J. (1993). Case-control studies with errors in covariates. *Journal of the American Statistical Association*, **88**, no. 421, 185–199.
- Carroll, R. J., Delaigle, A. & Hall, P. (2007). Non-parametric regression estimation from data contaminated by a mixture of berkson and classical errors. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **69**, no. 5, 859–878.

BIBLIOGRAPHY

- Carroll, R. J., Ruppert, D., Crainiceanu, C. M., Tosteson, T. D. & Karagas, M. R. (2004). Nonlinear and nonparametric regression and instrumental variables. *Journal of the American Statistical Association*, **99**, no. 467.
- Carroll, R. J., Ruppert, D., Stefanski, L. A. & Crainiceanu, C. M. (2006). *Measurement error in nonlinear models: a modern perspective*. CRC press.
- Cressie, N. (1993). *Statistics for spatial data: Wiley series in probability and statistics*.
- Diggle, P., Heagerty, P., Liang, K.-Y. & Zeger, S. (2002). *Analysis of longitudinal data*. Oxford University Press.
- Ezekiel, M. (1924). A method of handling curvilinear correlation for any number of variables. *Journal of the American Statistical Association*, **19**, no. 148, 431–453.
- Faes, C., Ormerod, J. & Wand, M. (2011). Variational bayesian inference for parametric and nonparametric regression with missing data. *Journal of the American Statistical Association*, **106**, no. 495.
- Fitzmaurice, G. M., Laird, N. M. & Ware, J. H. (2012). *Applied longitudinal analysis*, volume 998. John Wiley & Sons.
- Friedman, J. H. & Stuetzle, W. (1981). Projection pursuit regression. *Journal of the American statistical Association*, **76**, no. 376, 817–823.
- Ganguli, B., Staudenmayer, J. & Wand, M. (2005). Additive models with predictors subject to measurement error. *Australian & New Zealand Journal of Statistics*, **47**, no. 2, 193–202.
- Gelman, A. & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.
- Goldstein, H. (2011). *Multilevel statistical models*, volume 922. John Wiley & Sons.
- Green, P. & Silverman, B. (1994). Nonparametric regression and generalized linear models, vol. 58 of. *Monographs on Statistics and Applied Probability*.
- Hall, P., Pham, T., Wand, M. P., Wang, S. S. *et al.* (2011). Asymptotic normality and valid inference for gaussian variational approximation. *The Annals of Statistics*, **39**, no. 5, 2502–2532.

BIBLIOGRAPHY

- Hoffman, M., Bach, F. R. & Blei, D. M. (2010). Online learning for latent dirichlet allocation. In *advances in neural information processing systems*. pages 856–864.
- Huang, A., Wand, M. *et al.* (2013). Simple marginally noninformative prior distributions for covariance matrices. *Bayesian Analysis*, **8**, no. 2, 439–452.
- Jordan, M. I. (2004). Graphical models. *Statistical Science*, 140–155.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S. & Saul, L. K. (1999). An introduction to variational methods for graphical models. *Machine learning*, **37**, no. 2, 183–233.
- Knowles, D. A. & Minka, T. (2011). Non-conjugate variational message passing for multinomial and binary regression. In *Advances in Neural Information Processing Systems*. pages 1701–1709.
- Lauritzen, S. L. (1996). *Graphical models*. Oxford University Press.
- Lázaro-Gredilla, M. & Titsias, M. (2011). Variational heteroscedastic gaussian process regression. *Proceedings of the 28th International Conference on Machine Learning*.
- Lee, C. Y. & Wand, M. (2015). Streamlined mean field variational bayes for longitudinal and multilevel data analysis. *Unpublished manuscript*.
- Li, Y. & Ruppert, D. (2008). On the asymptotics of penalized splines. *Biometrika*, **95**, no. 2, 415–436.
- Liang, H., Wu, H. & Carroll, R. J. (2003). The relationship between virologic and immunologic responses in aids clinical research using mixed-effects varying-coefficient models with measurement error. *Biostatistics*, **4**, no. 2, 297–312.
- Ligges, U., Thomas, A., Spiegelhalter, D., Best, N. G., Lunn, D., Rice, K. & Strutz, S. (2009). *BRugs 0.5: OpenBUGS and Its R/S-PLUS Interface BRugs*. <http://www.stats.ox.ac.uk/pub/RWin/src/contrib/>.
- Lin, X. & Carroll, R. J. (2001). Semiparametric regression for clustered data using generalized estimating equations. *Journal of the American Statistical Association*, **96**, no. 455.
- Luenberger, D. G. & Ye, Y. (2008). *Linear and nonlinear programming*, volume 116. Springer.

BIBLIOGRAPHY

- Luts, J., Broderick, T. & Wand, M. P. (2013). Real-time semiparametric regression. *Journal of Computational and Graphical Statistics*, , no. just-accepted.
- Luts, J., Wang, S., Ormerod, J. & Wand, M. (2014). Semiparametric regression analysis via infer. net.
- Magnus, J. & Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics (revised ed.)*. New York: Wiley.
- Mallick, B., Hoffman, F. O. & Carroll, R. J. (2002). Semiparametric regression modeling with mixtures of berkson and classical error, with application to fallout from the nevada test site. *Biometrics*, **58**, no. 1, 13–20.
- Mardia, K. V., Kent, J. T. & Bibby, J. M. (1979). *Multivariate analysis*. Academic press.
- Marley, J. & Wand, M. (2010). Non-standard semiparametric regression via brugs. *Journal of Statistical Software*, **37**, no. 5.
- McGrory, C. A. & Titterton, D. (2007). Variational approximations in bayesian model selection for finite mixture distributions. *Computational Statistics & Data Analysis*, **51**, no. 11, 5352–5367.
- Minka, T. (2005). Divergence measures and message passing. Technical report, Technical report, Microsoft Research.
- Minka, T. & Winn, J. (2009). Gates. In D. Koller, D. Schuurmans, Y. Bengio & L. Bottou, editors, *Advances in Neural Information Processing Systems 21*. Curran Associates, Inc., pages 1073–1080.
- Minka, T., Winn, J., Guiver, G. & Kannan, A. (2008). *Infer.Net*. Microsoft Research Cambridge, Cambridge, UK.
- Müller, P. & Roeder, K. (1997). A bayesian semiparametric model for case-control studies with errors in variables. *Biometrika*, **84**, no. 3, 523–537.
- Neville, S. E., Ormerod, J. T. & Wand, M. P. (2014). Mean field variational bayes for continuous sparse signal shrinkage: Pitfalls and remedies. *Electronic Journal of Statistics*, **8**, no. 1, 1113–1151.

BIBLIOGRAPHY

- Nott, D. J., Tran, M.-N. & Leng, C. (2012). Variational approximation for heteroscedastic linear models and matching pursuit algorithms. *Statistics and Computing*, **22**, no. 2, 497–512.
- Opper, M. & Archambeau, C. (2009). The variational gaussian approximation revisited. *Neural computation*, **21**, no. 3, 786–792.
- Ormerod, J. & Wand, M. (2010). Explaining variational approximations. *The American Statistician*, **64**, no. 2, 140–153.
- Parisi, G. (1988). *Statistical field theory*. Addison-Wesley.
- Pebesma, E. J. (2004). Multivariate geostatistics in R: the `gstat` package. *Computers and Geosciences*, **30**, 683–691.
- Pebesma, E. J. & Duin, R. N. M. (2005). Spatial patterns of temporal change in north sea sediment quality on different spatial scales. in p. renard, h. demougeot-renard and r. froidevaux (eds). *Geostatistics for Environmental Application: Proceedings of the Fifth European Conference on Geostatistics for Environmental Applications (pp. 367-378)*.
- Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D. & The R Core Team (2009). *nlme 3.1: linear and nonlinear mixed effects models. R package. www.R-project.org*.
- Pratt, J. H., Jones, J. J., Miller, J. Z., Wagner, M. A. & Fineberg, N. S. (1989). Racial differences in aldosterone excretion and plasma aldosterone concentrations in children. *New England Journal of Medicine*, **321**, no. 17, 1152–1157. doi:10.1056/NEJM198910263211703. PMID: 2677724.
- R Core Team (2013). *R: a language and environment for statistical computing*. R foundation for Statistical Computing, Vienna, Austria, ISBN: 3-900051-07-0. www.R-project.org.
- Richardson, S., Leblond, L., Jaussent, I. & 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**, no. 3, 549–566.
- Rigby, R. & Stasinopoulos, D. (2005). Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **54**, no. 3, 507–554.

BIBLIOGRAPHY

- Robinson, G. K. (1991). That blup is a good thing: The estimation of random effects. *Statistical science*, 15–32.
- Roeder, K., Carroll, R. J. & Lindsay, B. G. (1996). A semiparametric mixture approach to case-control studies with errors in covariables. *Journal of the American Statistical Association*, **91**, no. 434, 722–732.
- Ruppert, D., Wand, M. & Carroll, R. J. (2009). Semiparametric regression during 2003–2007. *Electronic Journal of Statistics*, **3**, 1193.
- Ruppert, D., Wand, M. P. & Carroll, R. J. (2003). *Semiparametric regression*. New York: Cambridge University Press.
- Schafer, D. W. (2001). Semiparametric maximum likelihood for measurement error model regression. *Biometrics*, **57**, no. 1, 53–61.
- Simpson, D. (2013). Rodent tumor dataset monitoring frequency dependent backscatter coefficients. Department of Statistics, University of Illinois at Urbana-Champaign.
- Speed, T. (1991). [that blup is a good thing: The estimation of random effects]: Comment. *Statistical Science*, 42–44.
- Spiegelhalter, D. J., Thomas, A., Best, N. G., Gilks, W. R. & Lunn, D. (2003). *BUGS: Bayesian inference using Gibbs sampling*. Medical Research Council Biostatistics Unit, Cambridge, England. <http://www.mrc-bsu.cam.ac.uk/bugs>.
- Stan Development Team (2014). *Stan: A C++ Library for Probability and Sampling*. Version 2.2. <http://mc-stan.org>.
- Teschendorff, A. E., Wang, Y., Barbosa-Morais, N. L., Brenton, J. D. & Caldas, C. (2005). A variational bayesian mixture modelling framework for cluster analysis of gene-expression data. *Bioinformatics*, **21**, no. 13, 3025–3033.
- Thomas, A., O’Hara, B., Ligges, U. & Strutz, S. (2006). *Making BUGS Open*. R News. <http://CRAN.R-project.org/doc/Rnews/>.
- Titterton, D. (2004). Bayesian methods for neural networks and related models. *Statistical Science*, 128–139.
- Wahba, G. (1990). *Spline models for observational data*, volume 59. Siam.

BIBLIOGRAPHY

- Wainwright, M. J. & Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, **1**, no. 1-2, 1–305.
- Wand, M. (2009). Semiparametric regression and graphical models. *Australian & New Zealand Journal of Statistics*, **51**, no. 1, 9–41.
- Wand, M. (2015). Fast approximate inference for arbitrarily large semiparametric regression models via message passing.
- Wand, M. & Ormerod, J. (2008). On semiparametric regression with o’sullivan penalized splines. *Australian & New Zealand Journal of Statistics*, **50**, no. 2, 179–198.
- Wand, M. & Ripley, B. (2009). Kernsmooth: Functions for kernel smoothing for wand and jones (1995)kernel smoothing. r package.
- Wand, M. P. (2014). Fully simplified multivariate normal updates in non-conjugate variational message passing. *Journal of Machine Learning Research*, **15**, 1351–1369.
- Wand, M. P., Ormerod, J. T., Padoan, S. A., Fuhrwirth, R. *et al.* (2011). Mean field variational bayes for elaborate distributions. *Bayesian Analysis*, **6**, no. 4, 847–900.
- Wang, B. & Titterington, D. (2003). Local convergence of variational bayes estimators for mixing coefficients. *Preprint*.
- Wang, B., Titterington, D. *et al.* (2006). Convergence properties of a general algorithm for calculating variational bayesian estimates for a normal mixture model. *Bayesian Analysis*, **1**, no. 3, 625–650.
- Wang, C., Paisley, J. W. & Blei, D. M. (2011). Online variational inference for the hierarchical dirichlet process. In *International Conference on Artificial Intelligence and Statistics*. pages 752–760.
- Winn, J. M. & Bishop, C. M. (2005). Variational message passing. In *Journal of Machine Learning Research*. pages 661–694.