

Semiparametric Mean Field Variational Bayes: General Principles and Numerical Issues

David Rohde

*School of Mathematical and Physical Sciences
University of Technology Sydney
P.O. Box 123, Broadway, 2007, Australia*

DAVIDJROHDE@GMAIL.COM

Matt P. Wand

*School of Mathematical and Physical Sciences
University of Technology Sydney
P.O. Box 123, Broadway, 2007, Australia*

MATT.WAND@UTS.EDU.AU

Editor: Xiaotong Shen

Abstract

We introduce the term *semiparametric mean field variational Bayes* to describe the relaxation of mean field variational Bayes in which some density functions in the product density restriction are pre-specified to be members of convenient parametric families. This notion has appeared in various guises in the mean field variational Bayes literature during its history and we endeavor to unify this important topic. We lay down a general framework and explain how previous relevant methodologies fall within this framework. A major contribution is elucidation of numerical issues that impact semiparametric mean field variational Bayes in practice.

Keywords: Bayesian Computing, Factor Graph, Fixed-form Variational Bayes, Fixed-point Iteration, Non-conjugate Variational Message Passing, Nonlinear Conjugate Gradient Method

1. Introduction

We expound *semiparametric mean field variational Bayes*, a powerful combination of the notions of minimum Kullback-Leibler divergence and mean field restriction, that allows fast and often accurate approximate Bayesian inference for a wide range of scenarios. Most of its foundational literature and applications are in Machine Learning. However, semiparametric mean field variational Bayes is also an important paradigm for Statistics in the age of very big sample sizes and models.

Several articles concerned with deterministic approximate Bayesian inference, such as Barber and Bishop (1997), Honkela et al. (2010), Knowles and Minka (2011), Tan and Nott (2013), Wand (2014) and Menictas and Wand (2015), have demonstrated that modification of mean field variational Bayes (e.g. Wainwright and Jordan, 2008) to include pre-specified parametric families in the product density posterior approximation can have great practical benefits. For example, Barber and Bishop (1997) used a pre-specified Multivariate Normal distribution for the posterior density of the vector of adaptive parameters in multilayer neural networks while Tan and Nott (2013) derived a closed form variational approximate

algorithm for Bayesian Poisson mixed models by pre-specifying the fixed and random effects parameters to have Multivariate Normal distributions. Knowles and Minka (2011) took a message passing approach to mean field variational Bayes and explain how their approach to inclusion of pre-specified parametric families allows modular inference algorithms for arbitrary factors. Some recent articles on this topic have used the terms *fixed-form variational Bayes* (Honkela et al., 2010) and *non-conjugate variational message passing* (Knowles & Minka, 2011), to describe this modification of mean field variational Bayes. However, in this article we argue for adoption of the term *semiparametric mean field variational Bayes*.

Although we give a precise mathematical description of semiparametric mean field variational Bayes in Section 2, it simply refers to the relaxation of ordinary mean field variational Bayes in which some of the density functions in the postulated product density form are pre-specified to be particular parametric density functions, often chosen for reasons of tractability. This is a ‘halfway house’ between fully parametric approximation of the joint posterior density function of the model parameters with minimum Kullback-Leibler divergence used for parameter choice and pure mean field variational Bayes in which there is no parametric specification at all – only the product restriction. The following comments apply to our general framework:

- Semiparametric mean field variational Bayes is a modification of mean field variational Bayes that could be carried out via a message passing approach, as done by Knowles and Minka (2011), or by using the more common q -density approach used in, for example, Bishop (2006) and Ormerod and Wand (2010).
- The notion of conjugacy is not intrinsic to semiparametric mean field variational Bayes. The principle may be applied regardless of conjugacy relationships amongst the messages and/or q -densities. Therefore, the ‘non-conjugacy’ label used in recent articles for semiparametric relaxations of mean field variational Bayes is somewhat misleading.
- Contributions such as Knowles and Minka (2011) and Tan and Nott (2013) restrict attention to pre-specification of parametric families that are of exponential family form (e.g. Wainwright and Jordan, 2008). Whilst exponential family density functions have tractability advantages when used in semiparametric mean field variational Bayes, there is no intrinsic reason for only such densities to be used. In Section 5 we illustrate this point using pre-specified Skew-Normal density functions, which are not within the exponential family.
- Recent articles on non-conjugate variational message passing, such as Knowles and Minka (2011), Tan and Nott (2013) and Menictas and Wand (2015) used fixed-point iteration to minimize the Kullback-Leibler divergence or, equivalently, maximize the lower bound on the marginal log-likelihood. Theorem 1 of Knowles and Minka (2011) constitutes such an approach. However, any optimization approach could be used for obtaining the Kullback-Leibler optimal parameters such as Newton-Raphson iteration, quasi-Newton iteration, stochastic gradient descent, the Nelder-Mead simplex method and various hybrids and modifications of such methods.
- Some articles, such as the recent Challis and Barber (2013), are concerned solely with approximate inference via minimum divergence according to the pre-specification that

the posterior is within a parametric family such as Multivariate Normal with banded Cholesky covariance matrix factors. These contributions represent special cases of semiparametric mean field variational Bayes and their findings have relevance to the more general situation.

The main purposes of this article are:

- (1) Bring together the literature on semiparametric mean field variational Bayes and identify its core tenets.
- (2) Lay out and discuss numerical issues that arise in semiparametric mean field variational Bayes, which have a significant practical implications for this body of methodology.

The resulting exposition constitutes the first compendium on semiparametric mean field variational Bayes at its fullest level of generality. It can also be used as a basis for enhancements of semiparametric mean field variational Bayes methodology.

We use two examples to elucidate the general principles and numerical issues. The first, Example 1, involves a Bayesian model with a single parameter and, hence, is such that mean field approximation is not required. The simplicity of Example 1 allows a deep appreciation of the various issues with minimal notational overhead. Example 2 is the Bayesian Poisson mixed model treated in Wand (2014) and benefits from semiparametric mean field variational Bayes methodology. It demonstrates issues with high-dimensional optimization problems that are intrinsic to practical implementation.

One of the main outcomes of our numerical investigations is that fitting exponential family density functions via *natural* fixed-point iteration has some attractive properties. By ‘natural’, we mean a simple version of fixed-point iteration that arises when natural parametrizations are used. As we explain in Section 4, natural fixed-point iterations use *Riemannian* gradients to step through the parameter space, which is generally more efficient than ordinary gradients. The benefits of Riemannian gradient-based algorithms for Machine Learning problems goes back at least to Amari (1998). Such algorithms are the basis of the semiparametric mean field variational Bayes approach of Honkela et al. (2010).

In Section 2 we describe semiparametric mean field variational Bayes in full generality. A general overview of optimization strategies, pertinent to semiparametric mean field variational Bayes, is given in Section 3. The important special case of pre-specified exponential family density functions is treated in Section 4. Section 5 deals with the more difficult non-exponential family case via an illustrative example. Some closing remarks are given in Section 6.

2. General Principles

Semiparametric mean field variational Bayes is an approximate Bayesian inference method based on the principle of minimum Kullback-Leibler divergence. For arbitrary density functions p_1 and p_2 on \mathbb{R}^d ,

$$\text{KL}(p_1 \parallel p_2) \equiv \int_{\mathbb{R}^d} p_1(\mathbf{x}) \log \{p_1(\mathbf{x})/p_2(\mathbf{x})\} d\mathbf{x}$$

denotes the *Kullback-Leibler divergence* of p_2 from p_1 . Note that

$$\text{KL}(p_1 \parallel p_2) \geq 0 \quad \text{for any } p_1 \text{ and } p_2. \quad (1)$$

Consider a generic Bayesian model with observed data \mathcal{D} and parameter vector $(\boldsymbol{\theta}, \boldsymbol{\phi})$. The reason for this notational decomposition of the parameter vector will soon become apparent. Throughout this section we assume that $(\boldsymbol{\theta}, \boldsymbol{\phi})$ and \mathcal{D} are continuous random vectors with density functions $p(\boldsymbol{\theta}, \boldsymbol{\phi})$ and $p(\mathcal{D})$. The situation where some components are discrete has similar treatment with summations replacing integrals. Bayesian inference for $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ is based on the posterior density function

$$p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathcal{D}) = \frac{p(\mathcal{D}, \boldsymbol{\theta}, \boldsymbol{\phi})}{p(\mathcal{D})}.$$

The denominator, $p(\mathcal{D})$, is usually referred to as the *marginal likelihood* or the *model evidence*.

Let $q(\boldsymbol{\theta}, \boldsymbol{\phi})$ be an arbitrary density function over the parameter space of $(\boldsymbol{\theta}, \boldsymbol{\phi})$. The essence of variational approximate inference is to restrict $q(\boldsymbol{\theta}, \boldsymbol{\phi})$ to some class of density functions \mathcal{Q} and then use the optimal q -density function, given by

$$q^*(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL} \left\{ q(\boldsymbol{\theta}, \boldsymbol{\phi}) \parallel p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathcal{D}) \right\}, \quad (2)$$

as an approximation to the posterior density function $p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathcal{D})$.

Simple algebraic arguments (e.g. Section 2.1 of Ormerod and Wand, 2010) lead to

$$\log p(\mathcal{D}) = \text{KL} \left\{ q(\boldsymbol{\theta}, \boldsymbol{\phi}) \parallel p(\boldsymbol{\theta}, \boldsymbol{\phi} | \mathcal{D}) \right\} + \log \underline{p}(\mathcal{D}; q) \quad (3)$$

where

$$\underline{p}(\mathcal{D}; q) \equiv \exp \left[\int \int q(\boldsymbol{\theta}, \boldsymbol{\phi}) \log \left\{ \frac{p(\mathcal{D}, \boldsymbol{\theta}, \boldsymbol{\phi})}{q(\boldsymbol{\theta}, \boldsymbol{\phi})} \right\} d\boldsymbol{\theta} d\boldsymbol{\phi} \right]. \quad (4)$$

From (1) we have

$$\underline{p}(\mathcal{D}; q) \leq p(\mathcal{D}) \quad \text{for any } q(\boldsymbol{\theta}, \boldsymbol{\phi})$$

showing that $\underline{p}(\mathcal{D}; q)$ is a lower bound on the marginal likelihood. The non-negativity condition (1) means that an equivalent form for the optimal q -density function is

$$q^*(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underset{q \in \mathcal{Q}}{\text{argmax}} \underline{p}(\mathcal{D}; q). \quad (5)$$

This alternative optimization problem has the attractive interpretation of $q^*(\boldsymbol{\theta}, \boldsymbol{\phi})$ being chosen to maximize a lower bound on the marginal likelihood. For the remainder of this article we work with (5) rather than (2).

Parametric variational approximate inference involves setting

$$\mathcal{Q} = \{q(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\},$$

corresponding to a parametric family of density functions with parameter vector $\boldsymbol{\xi}$ ranging over Ξ . In this case (5) reduces to

$$q^*(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underset{\boldsymbol{\xi} \in \Xi}{\text{argmax}} \underline{p}(\mathcal{D}; q, \boldsymbol{\xi}), \quad (6)$$

where $\underline{p}(\mathcal{D}; q, \boldsymbol{\xi})$ is the marginal likelihood lower bound defined by (4), but with the dependence on $\boldsymbol{\xi}$ reflected in the notation.

An early contribution of this type is Hinton and van Camp (1993) who used minimum Kullback-Leibler divergence for Gaussian approximation of posterior density functions in neural networks models. Gaussian \mathcal{Q} families have also been used by Lappalainen and Honkela (2000), Archambeau et al. (2007), Raiko et al. (2007), Nickisch and Rasmussen (2008), Honkela and Valpola (2005), Honkela et al. (2007), Honkela et al. (2008) and Opper and Archambeau (2009). The recent contribution by Challis and Barber (2013) is an in-depth coverage of Gaussian minimum Kullback-Leibler approximate inference. Salimans and Knowles (2013) devised a stochastic approximation algorithm for solving (6) when \mathcal{Q} is a parametric family of exponential family form. Gershman et al. (2012) and Zobay (2014) investigated Gaussian-mixture extensions.

In what one may label a *nonparametric* variational approximation approach, ordinary mean field variational Bayes uses restricted q -density spaces such as

$$\mathcal{Q} = \{q(\boldsymbol{\theta}, \phi) : q(\boldsymbol{\theta}, \phi) = q(\boldsymbol{\theta}_1) \cdots q(\boldsymbol{\theta}_M) q(\phi)\} \text{ for some partition } \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\} \text{ of } \boldsymbol{\theta}. \quad (7)$$

The word ‘nonparametric’ is justified by the fact that there is no pre-specification that the q -density, or any of its factors, belong to a particular parametric family. Restriction of $q(\boldsymbol{\theta}, \phi)$ to a product density form is the only pre-specification being made. An iterative scheme for solving (5) under (7) follows from the last displayed equation given in Section 10.1.1 of Bishop (2006). The scheme is listed explicitly as Algorithm 1 of Ormerod and Wand (2010). Note that a simple adjustment that caters for $(\boldsymbol{\theta}, \phi)$, rather than $\boldsymbol{\theta}$, is required for notation being used here. Gershman et al. (2012) also use the word ‘nonparametric’ to describe a variational approximation approach. However, their methodology is parametric in the sense of the terminology that we are using here.

We propose that the term *semiparametric mean field variational Bayes* be used for restrictions of the form:

$$\mathcal{Q} = \{q(\boldsymbol{\theta}, \phi) : q(\boldsymbol{\theta}, \phi) = q(\boldsymbol{\theta}_1) \cdots q(\boldsymbol{\theta}_M) q(\phi; \boldsymbol{\xi}), \boldsymbol{\xi} \in \Xi\} \quad (8)$$

where $\{q(\phi; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\}$ is a pre-specified parametric family of density functions in ϕ . Under (8) there is no insistence on the $q(\boldsymbol{\theta}_i)$ having a particular parametric form. For models possessing particular conjugacy properties the optimal q -densities, $q^*(\boldsymbol{\theta}_i)$, will belong to relevant conjugate families. However, in general, the optimal q -densities of the $\boldsymbol{\theta}_i$ can assume arbitrary forms; see, for example, Figure 6 of Pham et al. (2013). The quality of a variational approximation is limited by the restrictions imposed by the particular choice of \mathcal{Q} . Semiparametric mean field variational Bayes imposes a product density restriction and then a parametric constraint on one of the factors. The overall quality of the approximation is determined by the combination of these two restrictions. While the estimated nonparametric factors are optimal given the product restriction, a parametric restriction with fewer product assumptions may be more accurate.

We now turn to the practical problem of solving the optimization problem (5) when the q -density restriction is of the form (8). Appropriate strategies for solving (5) depend on the nature of $q(\phi; \boldsymbol{\xi})$ as a function of $\boldsymbol{\xi}$ and the set Ξ . Some possibilities are:

- (A) Ξ is a finite set.

- (B) Ξ is an open subset of \mathbb{R}^d for some $d \in \mathbb{N}$ and $q(\phi; \xi)$ is smooth function of ξ over $\xi \in \Xi$.
- (C) Ξ is an open subset of \mathbb{R}^d for some $d \in \mathbb{N}$ and $q(\phi; \xi)$ is a non-smooth function of ξ over $\xi \in \Xi$.
- (D) Ξ is a complicated set that does not satisfy any of the descriptions given in (A)–(C).

For the vast majority of models in common use and $q(\phi; \xi)$ families (B) applies and most of the remainder of this article is devoted to that case. However, we will first briefly deal with (A) in Section 2.1, since it aids understanding of the semiparametric extension of mean field variational Bayes. To date, we are unaware of any semiparametric mean field variational Bayes contributions where (C) or (D) apply, so these cases are left aside.

2.1 Finite Parameter Space Case

Suppose that Ξ is a finite set. Then Algorithm 1 is the natural extension of the mean field variational Bayes coordinate ascent algorithm given, for example, in Section 10.1.1 of Bishop (2006) and Algorithm 1 of Ormerod and Wand (2010). In Algorithm 1, and elsewhere, the notation $\theta \setminus \theta_i$ denotes the vector θ with the entries of θ_i excluded.

For each $\xi \in \Xi$:

Initialize: $q(\theta_2), \dots, q(\theta_M)$.

Cycle:

$$q(\theta_1) \leftarrow \frac{\exp [E_{q(\theta \setminus \theta_1)} q(\phi; \xi) \{\log p(\mathbf{y}, \theta, \phi)\}]}{\int \exp [E_{q(\theta \setminus \theta_1)} q(\phi; \xi) \{\log p(\mathbf{y}, \theta, \phi)\}] d\theta_1}$$

⋮

$$q(\theta_M) \leftarrow \frac{\exp [E_{q(\theta \setminus \theta_M)} q(\phi; \xi) \{\log p(\mathbf{y}, \theta, \phi)\}]}{\int \exp [E_{q(\theta \setminus \theta_M)} q(\phi; \xi) \{\log p(\mathbf{y}, \theta, \phi)\}] d\theta_M}$$

until the increase in $\underline{p}(\mathcal{D}; q, \xi)$ is negligible.

$$q^*(\theta_i; \xi) \leftarrow q(\theta_i), \quad 1 \leq i \leq M \quad ; \quad \underline{p}(\mathcal{D}; q^*, \xi) \leftarrow \underline{p}(\mathcal{D}; q, \xi).$$

$$\xi^* \leftarrow \operatorname{argmax}_{\xi \in \Xi} \underline{p}(\mathcal{D}; q^*, \xi) \quad ; \quad q^*(\theta_i) \leftarrow q^*(\theta_i; \xi^*), \quad 1 \leq i \leq M.$$

Algorithm 1: *Coordinate ascent algorithm for semiparametric mean field variational Bayes when Ξ is a finite parameter space.*

For each value of ξ in Ξ , Algorithm 1 is essentially the ordinary mean field variational Bayes iterative algorithm — but with the density function of ϕ set to the parametric density function $q(\phi; \xi)$. The optimal ξ is then found by maximizing over the approximate marginal likelihood values that are recorded for each element of Ξ .

2.2 Infinite Parameter Space Case

Algorithm 1 shows how to obtain the Kullback-Leibler-optimal $q(\theta_i)$ and $q(\phi; \boldsymbol{\xi})$ densities in the case where Ξ is finite. However, for common parametric families such as the Normal and Gamma, Ξ is infinite and the solution of (5) under (8) is more delicate. The coordinate ascent idea used to obtain the $q^*(\theta_i)$ in Algorithm 1 can still be entertained. However, it needs to be combined with an optimization scheme that searches for the optimal $\boldsymbol{\xi}$ over the infinite space Ξ .

For the remainder of this article we focus on the problem of solving (5) under restriction (B) on the q -density parameter space Ξ . We start by studying the criterion function $p(\mathcal{D}; q, \boldsymbol{\xi})$ and special forms that it takes under the mean field restriction. The notions of entropy and factor graphs are shown to be very relevant and useful. We then introduce two running examples, Example 1 and Example 2, to illustrate the issues involved. Since Example 1 has only one parameter requiring inference, this is not a fully-fledged semiparametric mean field variational Bayes problem and the optimization problem is of the form (6). Additionally, (6) for Example 1 is a bivariate optimization problem which allows deeper probing of the numerical analytic issues. Example 2 uses the Poisson mixed model, treated in Section 5.1 of Wand (2014), as our main semiparametric mean field variational Bayes example. It is substantial enough to convey various practical issues but also has a closed form $\log p(q; \boldsymbol{\xi})$ expression that allows purely algebraic exposition.

2.2.1 ENTROPY, FACTOR GRAPHS AND THE MARGINAL LOG-LIKELIHOOD LOWER BOUND

If \boldsymbol{x} is a random vector having density function p then the corresponding *entropy* is given by

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

For many common distributional families, the entropy admits an algebraic expression in terms of the distribution's parameters. For example, if

$$p(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\}$$

is the Multivariate Normal density function of dimension d with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ then

$$\text{Entropy}(p; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2} d \{1 + \log(2\pi)\} + \frac{1}{2} \log |\boldsymbol{\Sigma}|. \quad (9)$$

Note that the entropy is independent of the mean vector $\boldsymbol{\mu}$.

Another entropy expression, which arises in many Bayesian models and the example of Section 2.2.3, is that for the Inverse Gamma family of density functions. Let

$$x \sim \text{Inverse-Gamma}(\kappa, \lambda)$$

denote the random variable x having density function

$$p(x; \kappa, \lambda) = \frac{\kappa^\lambda}{\Gamma(\kappa)} x^{-\kappa-1} \exp(-\lambda/x), \quad x > 0,$$

with parameters $\kappa, \lambda > 0$. In this case

$$\text{Entropy}(p; \kappa, \lambda) = \log(\lambda) + \kappa + \log\{\Gamma(\kappa)\} - (\kappa + 1)\text{digamma}(\kappa) \quad (10)$$

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

The next relevant concept is that of a *factor graph*, which we first explain via an example. Consider the approximate Bayesian inference problem according to the semiparametric mean field variational Bayes restriction (8). Figure 1 is the factor graph for an $M = 9$ example of (8) with the joint density function of all random vectors in the model factorizing as follows:

$$p(\mathbf{x}, \boldsymbol{\theta}, \phi) = f_1(\boldsymbol{\theta}_1) f_2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \phi) f_3(\boldsymbol{\theta}_3) f_4(\boldsymbol{\theta}_4) f_5(\boldsymbol{\theta}_5, \phi) f_6(\boldsymbol{\theta}_5) f_7(\boldsymbol{\theta}_6, \phi) f_8(\boldsymbol{\theta}_6) \\ \times f_9(\boldsymbol{\theta}_7, \boldsymbol{\theta}_8, \boldsymbol{\theta}_9) f_{10}(\boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \boldsymbol{\theta}_4) \quad (11)$$

for *factors* f_1, \dots, f_{10} . Note that some of these factors depend on the data vector \mathbf{x} , but the dependence is suppressed in the f_j notation. Specific examples are given in Sections 2.2.2 and 2.2.3. In Figure 1 the circles correspond to the components of the mean field product restriction

$$q(\phi) \prod_{i=1}^9 q(\boldsymbol{\theta}_i) \quad (12)$$

and are called *stochastic nodes*. The solid squares correspond to the factors f_j , $1 \leq j \leq 10$, and are called *factor nodes*. Edges join each factor node f_j to those stochastic nodes that are included in the f_j function.

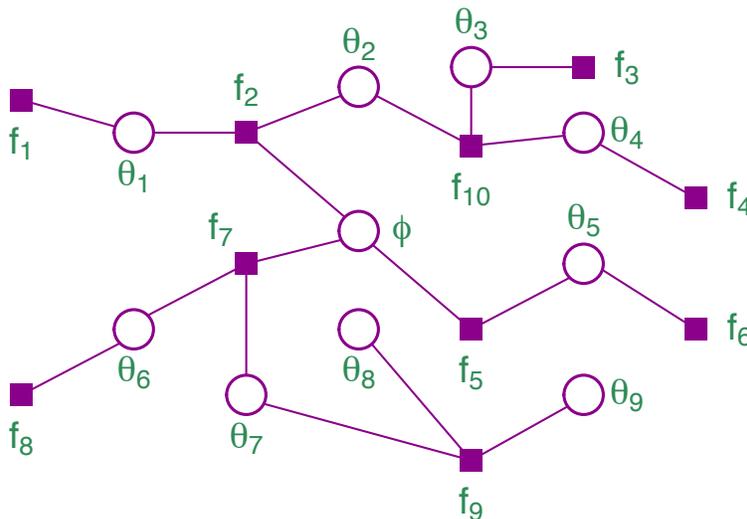


Figure 1: *Factor graph corresponding to the model (11) and q-density product restriction (12).*

Now consider the general case with semiparametric mean field restriction (8) and suppose that $p(\mathbf{x}, \boldsymbol{\theta}, \phi)$ has N factors f_j , $1 \leq j \leq N$. Then the marginal log-likelihood lower bound admits the following expression in terms of the components of the corresponding factor graph:

$$\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi}) = \text{Entropy}\{q(\phi; \boldsymbol{\xi})\} + \sum_{i=1}^M \text{Entropy}\{q(\boldsymbol{\theta}_i)\} + \sum_{j=1}^N E_q\{\log(f_j)\}.$$

The ϕ -localized component of $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})$, which we denote by $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]}$, is

$$\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]} \equiv \text{Entropy}\{q(\phi; \boldsymbol{\xi})\} + \sum_{j \in \text{neighbors}(\phi)} E_q\{\log(f_j)\} \quad (13)$$

where

$$\begin{aligned} \text{neighbors}(\phi) &\equiv \{1 \leq j \leq N : f_j \text{ is a neighbor of } \phi \text{ on the factor graph}\} \\ &= \{1 \leq j \leq N : f_j \text{ involves } \phi\}. \end{aligned}$$

For the factor graph shown in Figure 1 $\text{neighbors}(\phi) = \{2, 5, 7\}$ and so have

$$\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]} = \text{Entropy}\{q(\phi; \boldsymbol{\xi})\} + E_q\{\log(f_2)\} + E_q\{\log(f_5)\} + E_q\{\log(f_7)\}.$$

as the ϕ -localized component of $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})$.

For large Bayesian models it is prudent to maximize this localized approximate log-likelihood as part of a coordinate ascent scheme involving all q -density parameters. Such an approach, combined with the *locality property* of mean field variational Bayes (e.g. Wand et al., 2011, Section 3), allows for streamlined handling of arbitrarily large models. We formalize this approach to semiparametric mean field variational Bayes in the shape of Algorithm 2 in the upcoming Section 2.2.4. However, we first give some concrete examples of mean field variational Bayes with pre-specified parametric family q -density functions.

2.2.2 EXAMPLE 1: GUMBEL RANDOM SAMPLE

A Bayesian model for a random sample x_1, \dots, x_n from a Gumbel distribution with location parameter ϕ and unit scale is

$$p(x_1, \dots, x_n | \phi) = \prod_{i=1}^n \exp\{-(x_i - \phi) - e^{-(x_i - \phi)}\}, \quad \phi \sim N(\mu_\phi, \sigma_\phi^2). \quad (14)$$

The posterior density function of ϕ is

$$p(\phi | \mathbf{x}) = \frac{\exp\left\{n\phi - e^\phi \sum_{i=1}^n e^{-x_i} - \frac{1}{2\sigma_\phi^2}(\phi - \mu_\phi)^2\right\}}{\int_{-\infty}^{\infty} \exp\left\{n\phi' - e^{\phi'} \sum_{i=1}^n e^{-x_i} - \frac{1}{2\sigma_\phi^2}(\phi' - \mu_\phi)^2\right\} d\phi'}. \quad (15)$$

The denominator on the right-hand side of (15) is not available in closed form. This implies that numerical methods such as quadrature are required to obtain the Bayes estimate of ϕ

and corresponding credible sets. Instead we consider minimum Kullback-Leibler divergence approximation of $p(\phi|\mathbf{x})$ over a parametric pre-specified family. Let

$$\mathcal{Q} = \{q(\phi; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\}$$

be such a parametric family. Then the optimal q -density is $q(\phi; \boldsymbol{\xi}^*)$ where

$$\boldsymbol{\xi}^* = \operatorname{argmax}_{\boldsymbol{\xi} \in \Xi} \underline{p}(q; \boldsymbol{\xi}). \quad (16)$$

Figure 2 shows the factor graph of the model, with factors $p(\phi)$ and $p(\mathbf{x}|\phi)$ neighboring the stochastic node ϕ .

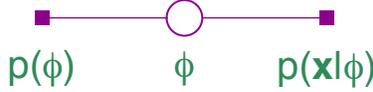


Figure 2: *Factor graph for the Example 1 model.*

The marginal log-likelihood lower bound is

$$\log \underline{p}(\mathbf{x}; \boldsymbol{\xi}) = \text{Entropy}\{q(\phi; \boldsymbol{\xi})\} + E_q\{\log p(\phi)\} + E_q\{\log p(\mathbf{x}|\phi)\} \quad (17)$$

and the contributions to $\log \underline{p}(\mathbf{x}; \boldsymbol{\xi})$ from the factors are

$$\begin{aligned} E_q\{\log p(\phi)\} &= -\frac{1}{2\sigma_\phi^2} \left[\{E_{q(\phi; \boldsymbol{\xi})}(\phi) - \mu_\phi\}^2 + \text{Var}_{q(\phi; \boldsymbol{\xi})}(\phi) \right] \\ \text{and } E_q\{\log p(\mathbf{x}|\phi)\} &= n E_{q(\phi; \boldsymbol{\xi})}(\phi) - M_{q(\phi; \boldsymbol{\xi})}(1) \sum_{i=1}^n e^{-x_i} - n\bar{x} \end{aligned} \quad (18)$$

where $M_{q(\phi; \boldsymbol{\xi})}$ is the moment generating function corresponding to $q(\phi; \boldsymbol{\xi})$.

Now suppose that

$$\mathcal{Q} = \left\{ q(\phi; \mu_{q(\phi)}, \sigma_{q(\phi)}^2) = \frac{1}{\sqrt{2\pi\sigma_{q(\phi)}^2}} \exp \left\{ \frac{-(\phi - \mu_{q(\phi)})^2}{2\sigma_{q(\phi)}^2} \right\} : \mu_{q(\phi)} \in \mathbb{R}, \sigma_{q(\phi)}^2 > 0 \right\}$$

corresponding to the Normal family with q -density parameter vector $\boldsymbol{\xi} = (\mu_{q(\phi)}, \sigma_{q(\phi)}^2)$ and parameter space $\Xi = \mathbb{R} \times \mathbb{R}_+$ where $\mathbb{R}_+ \equiv (0, \infty)$ is the positive half-line. Then, from the entropy result (9) and the well-known expression for the moment generating function of the Normal distribution we obtain

$$\begin{aligned} \log \underline{p}(\mathbf{x}; \mu_{q(\phi)}, \sigma_{q(\phi)}^2) &= \frac{1}{2} \{1 + \log(2\pi)\} + \frac{1}{2} \log(\sigma_{q(\phi)}^2) + n \mu_{q(\phi)} \\ &\quad - \exp(\mu_{q(\phi)} + \frac{1}{2}\sigma_{q(\phi)}^2) \sum_{i=1}^n e^{-x_i} - \frac{1}{2\sigma_\phi^2} \{(\mu_{q(\phi)} - \mu_\phi)^2 + \sigma_{q(\phi)}^2\} - n\bar{x}. \end{aligned}$$

It follows that the Kullback-Leibler optimal Normal q -density function is $q(\phi; \mu_{q(\phi)}^*, (\sigma_{q(\phi)}^2)^*)$ where

$$\begin{bmatrix} \mu_{q(\phi)}^* \\ (\sigma_{q(\phi)}^2)^* \end{bmatrix} = \operatorname{argmax}_{\mu_{q(\phi)} \in \mathbb{R}, \sigma_{q(\phi)}^2 > 0} \left\{ f_{\text{Ex1}}^N \left(\mu_{q(\phi)}, \sigma_{q(\phi)}^2; n, \sum_{i=1}^n e^{-x_i}, \mu_\phi, \sigma_\phi^2 \right) \right\} \quad (19)$$

and

$$f_{\text{Ex1}}^N(x, y; a, b, c, d) = \frac{1}{2} \log(y) + ax - b \exp(x + \frac{1}{2}y) - \frac{1}{2d} \{(x - c)^2 + y\}. \quad (20)$$

The main arguments satisfy $x \in \mathbb{R}$, $y > 0$ and the auxiliary arguments are such that $a, b, d > 0$ and $c \in \mathbb{R}$. From (19) we see that the minimum Kullback-Leibler divergence problem (16), where \mathcal{Q} is the Normal family, reduces to a non-linear bivariate optimization problem. Theory given in Challis and Barber (2013) applies to this example. For example, results given in their Section 3.2 can be used to establish that $f_{\text{Ex1}}^N(x, y; a, b, c, d)$ is jointly concave in x and \sqrt{y} .

In Section 3 we study strategies for solving such problems and apply them to this example in Section 4.2.

2.2.3 EXAMPLE 2: POISSON MIXED MODEL

A single variance component *Poisson mixed model* is

$$\begin{aligned} y_i | \boldsymbol{\beta}, \mathbf{u} &\overset{\text{ind.}}{\sim} \text{Poisson}[\exp\{(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_i\}], \quad 1 \leq i \leq n, \\ \mathbf{u} | \sigma^2 &\sim N(0, \sigma^2 \mathbf{I}_K), \quad \sigma^2 | a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/a), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}_p), \quad a \sim \text{Inverse-Gamma}(\frac{1}{2}, 1/A^2) \end{aligned} \quad (21)$$

where \mathbf{X} is an $n \times p$ fixed effects design matrix, \mathbf{Z} is an $n \times K$ random effects design matrix. Note that the prior on σ in (21) is the Half Cauchy distribution with scale parameter A :

$$p(\sigma) = \frac{(2/\pi)}{A\{1 + (\sigma/A)^2\}}, \quad \sigma > 0.$$

In (21) $\sigma_{\boldsymbol{\beta}} > 0$ and $A > 0$ are hyperparameters to be chosen by the analyst.

A mean field approximation to the joint posterior density function of the model parameters is

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, a | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma^2) q(a). \quad (22)$$

As detailed in Appendix A.3, the optimal q -density functions satisfy

$$\begin{aligned} q^*(\sigma^2) \text{ and } q^*(a) &\text{ are both Inverse-Gamma density functions, and} \\ q^*(\boldsymbol{\beta}, \mathbf{u}) &\propto \exp \left\{ \mathbf{y}^T (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) - \mathbf{1}^T \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \right. \\ &\quad \left. - \frac{1}{2} \sigma_{\boldsymbol{\beta}}^{-2} \|\boldsymbol{\beta}\|^2 - \frac{1}{2} E_{q(\sigma^2)}(1/\sigma^2) \|\mathbf{u}\|^2 \right\}. \end{aligned} \quad (23)$$

Since $q^*(\boldsymbol{\beta}, \mathbf{u})$ is not a standard form, numerical methods are required to obtain the variational approximate Bayes estimates and credible sets. *Semiparametric* mean field variational Bayes alternatives take the form

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, a | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) q(\sigma^2) q(a) \quad (24)$$

where $\{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\}$ is a pre-specified parametric family of density functions. The optimal density functions $q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}^*)$, $q^*(\sigma^2)$ and $q^*(a)$ are found by minimizing

$$\text{KL} \left\{ q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) q(\sigma^2) q(a) \parallel p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, a | \mathbf{y}) \right\}. \quad (25)$$

We now focus on solving (25).

In Appendix C of Wand (2014) it is shown that

$q^*(\sigma^2)$ is an Inverse-Gamma($\frac{1}{2}(K+1), B_{q(\sigma^2)}$) density function, and
 $q^*(a)$ is an Inverse-Gamma($1, B_{q(a)}$) density function

where

$$B_{q(\sigma^2)} = \frac{1}{2} [\| E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u}) \|^2 + \text{tr}\{\text{Cov}_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u})\}] + \mu_{q(1/a)} \quad (26)$$

and $B_{q(a)} = \mu_{q(1/\sigma^2)} + A^{-2}$ with the definition

$$\mu_{q(1/v)} \equiv E_{q(v)}(1/v)$$

for a generic random variable v .

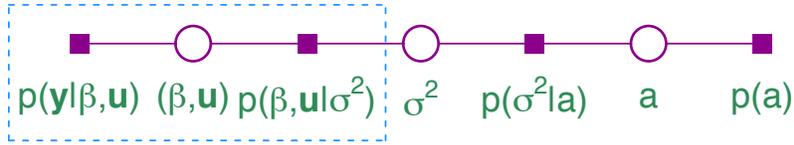


Figure 3: *Factor graph for the Example 2 model with stochastic nodes corresponding the mean field restriction (24). The dashed line box contains the stochastic node (β, \mathbf{u}) and its neighboring factors.*

It remains to obtain the optimal value of $\boldsymbol{\xi}$ in $q(\beta, \mathbf{u}; \boldsymbol{\xi})$. Figure 3 shows the factor graph of the current model under mean field restriction (24). The lower bound on the marginal log-likelihood, in terms of the stochastic nodes and factors of Figure 3, is

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi}) &= \text{Entropy}\{q(\beta, \mathbf{u}; \boldsymbol{\xi})\} + \text{Entropy}\{q(\sigma^2)\} + \text{Entropy}\{q(a)\} \\ &\quad + E_q\{\log p(\mathbf{y}|\beta, \mathbf{u})\} + E_q\{\log p(\beta, \mathbf{u}|\sigma^2)\} \\ &\quad + E_q\{\log p(\sigma^2|a)\} + E_q\{\log p(a)\}. \end{aligned} \quad (27)$$

One could substitute (26) into the $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi})$ expression. This resulting marginal log-likelihood lower bound would depend on the q -densities only through $\boldsymbol{\xi}$ and could then be maximized over $\boldsymbol{\xi} \in \Xi$.

An alternative strategy, that scales better to larger models, is to use a coordinate ascent scheme that maximizes $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi})^{[(\beta, \mathbf{u})]}$, the (β, \mathbf{u}) -localized component of $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi})$, over $\boldsymbol{\xi} \in \Xi$. The relevant factors are those neighboring (β, \mathbf{u}) in Figure 3, corresponding to the dashed line box. The quantity requiring maximization is then

$$\begin{aligned} \log \underline{p}(\mathbf{y}; q, \boldsymbol{\xi})^{[(\beta, \mathbf{u})]} &= \text{Entropy}\{q(\beta, \mathbf{u}; \boldsymbol{\xi})\} + E_q\{\log p(\mathbf{y}|\beta, \mathbf{u})\} + E_q\{\log p(\beta, \mathbf{u}|\sigma^2)\} \\ &= \text{Entropy}\{q(\beta, \mathbf{u}; \boldsymbol{\xi})\} + \mathbf{y}^T \{\mathbf{X} E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\beta) + \mathbf{Z} E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u})\} \\ &\quad - \mathbf{1}^T E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})} \{\exp(\mathbf{X}\beta + \mathbf{Z}\mathbf{u})\} \\ &\quad - \frac{1}{2\sigma_\beta^2} [\| E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\beta) \|^2 + \text{tr}\{\text{Cov}_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\beta)\}] \\ &\quad - \frac{1}{2} \mu_{q(1/\sigma^2)} [\| E_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u}) \|^2 + \text{tr}\{\text{Cov}_{q(\beta, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u})\}] + \text{const} \end{aligned} \quad (28)$$

where ‘const’ denotes terms not depending on ξ .

Next, suppose that \mathcal{Q} corresponds to the family of Multivariate Normal density functions in (β, \mathbf{u}) :

$$q(\beta, \mathbf{u}; \boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}) = (2\pi)^{-(p+K)/2} |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}|^{-1/2} \\ \times \exp \left\{ -\frac{1}{2} \left(\begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} - \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \right)^T \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}^{-1} \left(\begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} - \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \right) \right\}.$$

Then the (β, \mathbf{u}) -localized approximate marginal log-likelihood reduces to

$$\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})^{[(\beta, \mathbf{u})]} = \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})}| + \mathbf{y}^T \mathbf{C} \boldsymbol{\mu}_{q(\beta, \mathbf{u})} \\ - \mathbf{1}^T \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\beta)} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} \mathbf{C}^T) \right\} \\ - \frac{1}{2\sigma_\beta^2} \left\{ \|\boldsymbol{\mu}_{q(\beta)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\beta)}) \right\} \\ - \frac{1}{2} \mu_{q(1/\sigma^2)} \left\{ \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})}) \right\} + \text{const} \quad (29)$$

where $\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$, $\text{diagonal}(\mathbf{M})$ is the vector containing the diagonal entries of the square matrix \mathbf{M} , and $\boldsymbol{\mu}_{q(\beta)}$ is the sub-vector of $\boldsymbol{\mu}_{q(\beta, \mathbf{u})}$ corresponding to β . Analogous definitions apply to $\boldsymbol{\mu}_{q(\mathbf{u})}$, $\boldsymbol{\Sigma}_{q(\beta)}$ and $\boldsymbol{\Sigma}_{q(\mathbf{u})}$. Appendix A.3 provides derivational details for (29).

For this example, the full coordinate ascent scheme has updates as follows:

perform one or more updates of $(\boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})$ within an iterative scheme for the optimization problem:

$$\underset{\boldsymbol{\mu}'_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}'_{q(\beta, \mathbf{u})}}{\text{argmax}} \left\{ \log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}'_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}'_{q(\beta, \mathbf{u})})^{[(\beta, \mathbf{u})]} \right\} \\ B_{q(a)} \leftarrow \mu_{q(1/\sigma^2)} + A^{-2} \quad ; \quad \mu_{q(1/a)} \leftarrow 1/B_{q(a)} \\ B_{q(\sigma^2)} \leftarrow \frac{1}{2} [\|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}\{\boldsymbol{\Sigma}_{q(\mathbf{u})}\}] + \mu_{q(1/a)} \\ \mu_{q(1/\sigma^2)} \leftarrow \frac{1}{2}(K+1)/B_{q(\sigma^2)}.$$

For now, we deliberately leave the form of the $(\boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})$ maximization strategy unspecified. We also allow for one or more updates of an iterative scheme aimed at maximizing $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}'_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}'_{q(\beta, \mathbf{u})})^{[(\beta, \mathbf{u})]}$, rather than iterating to convergence at every iteration of the full coordinate ascent scheme. Section 3 describes various optimization strategies that could be used for updating $(\boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})$.

Negligible absolute change in $\log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})$ can be used as a stopping criterion for the iterations and an algebraic expression for this quantity is given in Appendix A.3.

We return to Example 2 in Section 4.3.

2.2.4 GENERAL SEMIPARAMETRIC MEAN FIELD VARIATIONAL BAYES ALGORITHM

We now treat semiparametric mean field variational Bayes in general, with the set-up laid out in Section 2 and restriction (8). Let $\log \underline{p}(\mathcal{D}; q, \xi)^{[\phi]}$ be defined with respect to the

factor graph of $p(\mathbf{x}, \phi, \boldsymbol{\theta})$ with stochastic nodes $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M$ and ϕ . Algorithm 2 is a general coordinate ascent algorithm for approximate inference that builds on the standard mean field variational Bayes algorithm.

Initialize: $q(\boldsymbol{\theta}_2), \dots, q(\boldsymbol{\theta}_M)$.

Cycle:

perform one or more updates of $\boldsymbol{\xi}$ within an iterative scheme for the optimization problem:

$$\begin{aligned} & \operatorname{argmax}_{\boldsymbol{\xi}' \in \Xi} \left\{ \log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi}')^{[\phi]} \right\} \\ q(\boldsymbol{\theta}_1) & \leftarrow \frac{\exp \left[E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_1) q(\phi; \boldsymbol{\xi})} \{ \log p(\mathbf{y}, \boldsymbol{\theta}, \phi) \} \right]}{\int \exp \left[E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_1) q(\phi; \boldsymbol{\xi})} \{ \log p(\mathbf{y}, \boldsymbol{\theta}, \phi) \} \right] d\boldsymbol{\theta}_1} \\ & \quad \vdots \\ q(\boldsymbol{\theta}_M) & \leftarrow \frac{\exp \left[E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_M) q(\phi; \boldsymbol{\xi})} \{ \log p(\mathbf{y}, \boldsymbol{\theta}, \phi) \} \right]}{\int \exp \left[E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_M) q(\phi; \boldsymbol{\xi})} \{ \log p(\mathbf{y}, \boldsymbol{\theta}, \phi) \} \right] d\boldsymbol{\theta}_M} \end{aligned}$$

until the absolute change in $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})$ is negligible.

Algorithm 2: *The general semiparametric mean field variational Bayes algorithm for restriction (8) with $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]}$ defined with respect to factor graph of $p(\mathbf{x}, \boldsymbol{\theta}, \phi)$ with stochastic nodes $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M$ and ϕ .*

For each of these approaches, there remains the practical problem of devising an optimization scheme for $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]}$ and ensuring that it leads to the optimal parameters being chosen. Section 3 deals with this problem.

2.3 Relationship to Existing Literature

The general principle of semiparametric mean field variational Bayes that we have laid out in this section is not brand new and, in fact, instances of this principle have made appearances in the literature since the late 1990s — although they are few in number. We now briefly survey articles known to us that have a semiparametric mean field variational Bayes component. As we will see, the terminology varies quite considerably.

Barber and Bishop (1997) uses the terms *ensemble learning* and *hyperparameter adaptation* for what essentially is a semiparametric mean field variational Bayes approach to fitting multi-layer neural networks. They pre-specify Multivariate Normal distributions for the coefficient vector but, in their Section 2.1, allow covariance matrix parameters to be unspecified except for mean field assumptions. However, they do not include numerical details for minimizing Kullback-Leibler divergence over the Multivariate Normal parameters.

Honkela et al. (2010) adopt the phrase *fixed-form variational Bayes* in what is a quite general approach to semiparametric mean field variational Bayes as summarized in their Algorithm 1. Optimization of the pre-specified parametric component parameters is achieved using the *nonlinear conjugate gradient method*, which we describe in Section 3.2. However, Honkela et al. (2010) work with *Riemannian* gradients which, they argue, are more efficient than Euclidean gradients.

In Knowles and Minka (2011) the focus is incorporation of pre-specified exponential family distributions whilst preserving the message passing aspect of a modular approach to mean field variational Bayes, known as *variational message passing*. They arrive at an extension which they label *non-conjugate variational message passing*. The exponential family distribution parameters are chosen via fixed-point iteration, which we describe in detail in Section 3.1.

Tan and Nott (2013) take a semiparametric mean field variational Bayes approach to approximate inference in Bayesian generalized linear mixed models for grouped data. They use pre-specified Multivariate Normal density functions for the random effects of each group with mean field product restrictions and achieve good approximation accuracy via so-called partially noncentered parameterizations.

In the case of pre-specified Multivariate Normal density functions, Wand (2014) obtains an explicit form for the fixed-point iteration scheme of Knowles and Minka (2011) and illustrated its use for the Poisson mixed model described in Section 2.2.3. In Luts and Wand (2015) and Menictas and Wand (2015), semiparametric mean field variational Bayes with Multivariate Normal pre-specification is applied, respectively, to count response semiparametric regression and heteroscedastic semiparametric regression.

2.4 Advantages and Limitations of Algorithm 2

As just described in Section 2.3, Algorithm 2 is a very useful generalization of the ordinary mean field variational Bayes algorithm and allows for tractable handling of a wider class of models. For example, the heteroscedastic nonparametric regression model of Menictas and Wand (2015) is such that ordinary mean field variational Bayes is numerically challenging if one uses the same product restrictions as used in homoscedastic nonparametric regression. The semiparametric mean field variational Bayes extension, based on Multivariate Normal pre-specification of basis function coefficients, leads to an iterative scheme with closed form updates. Simulation studies show very good accuracy compared with Markov chain Monte Carlo-based inference. However Algorithm 2 is not guaranteed to converge and, when it does converge, may result in mediocre approximate Bayesian inference. We close this section by briefly discussing such limitations of Algorithm 2.

In cases where a generic iterative procedure is used to solve $\operatorname{argmax}_{\xi' \in \Xi} \{\log p(\mathcal{D}; q, \xi')\}$ there is no guarantee that the lower bound is increased in a particular iteration or of convergence in general. As a consequence, the convergence guarantees enjoyed by ordinary mean field variational Bayes algorithms are not shared by their semiparametric extension. As we demonstrate in Section 4.2, convergence does not occur for particular numerical optimization strategies.

Lastly, there is the limitation imposed by the mean field restriction. Even though mean field variational Bayes can lead to very accurate approximate inference, there are

circumstances where its accuracy is quite poor. Some examples, with explanations for the inaccuracy, are given in Wang and Titterton (2005) and Neville et al. (2014). Semiparametric mean field variational Bayes shares this limitation since parametric pre-specification imposes a degradation in accuracy compared with ordinary mean field variational Bayes.

3. Numerical Optimization Strategies

In ordinary mean field variational Bayes, parameter optimization is achieved using a convex optimization algorithm that converges under reasonable assumptions (e.g. Luenberger and Ye, 2008). In the semiparametric extension there is no such convex optimization theory and general numerical optimization has to be called upon to optimize the parameters in the pre-specified parametric density function.

Numerical optimization is a major area of mathematical study with an enormous literature. Recent summaries are given in, for example, Givens and Hoeting (2005) and Ackleh et al. (2010), with the former being geared towards optimization problems arising in Statistics. The choice of optimization method typically is driven by factors such as the smoothness of the function requiring optimization and availability of expressions for low-order derivatives. Optimization methods with derivative information invariably take the form of iterative schemes. Semiparametric mean field variational Bayes often has the luxuries of smoothness and derivative expressions. It is also beneficial to have relatively simple iterative updates given the overarching goal of fast approximate inference. Therefore, we gear our summary of numerical optimization strategies towards simple derivative-based schemes. This is in keeping with the existing literature on parametric and semiparametric variational inference.

Let $f : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$ be a function and consider the problem of finding its maximum or minimum value over a set $D_0 \subseteq D$. If all partial derivatives of f exist then a necessary condition for a maximum or minimum in the interior of D_0 is the *stationary point condition*

$$\frac{\partial}{\partial x_j} f(\mathbf{x}) = 0, \quad 1 \leq j \leq d. \quad (30)$$

This converts the optimization problem to a multivariate root-finding problem. However (30) is not a sufficient condition for global optima since local optima and saddle points of f inside D_0 also satisfy (30). Properties of f , such as regions where it is concave or convex, can aid the determination of global optima.

Throughout this section we make use of the derivative matrix and Hessian matrix notation defined in Appendix A.2.

3.1 Fixed-Point Iteration

Assuming that the derivative vector $D_{\mathbf{x}}f(\mathbf{x})$ (defined formally in Appendix A.1) exists, the stationary point condition (30) can be written

$$D_{\mathbf{x}}f(\mathbf{x})^T = \mathbf{0} \quad (31)$$

where $\mathbf{0}$ is the vector of zeroes. Fixed-point iteration aims to find points that satisfy (31), which we denote generically by \mathbf{x}^* . Such points are then candidates as maxima or minima

of f . Firstly, (31) is rewritten in the form

$$\mathbf{x} = \mathbf{g}(\mathbf{x}) \tag{32}$$

for some function $\mathbf{g} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^d$. Given this \mathbf{g} , fixed-point iteration simply involves repeated evaluation of \mathbf{g} , as given in Algorithm 3

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

Cycle:

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

until convergence.

Algorithm 3: *The fixed-point iteration algorithm in generic form.*

Note, however, the following issues regarding fixed point iteration:

- For a given stationary point condition (31) there are numerous functions \mathbf{g} for which (32) holds. In other words, there are many possible fixed point algorithms available to solve (31).
- Once \mathbf{g} and \mathbf{x}_{init} are chosen then the above algorithm is *not necessarily guaranteed* to converge to a stationary point \mathbf{x}^* . There is a large literature on convergence of fixed-point iterative algorithms and good references on the topic include Section 8.1 of Ortega (1990) and Section 8.2 of Ackleh et al. (2010). For example Theorem 8.1.7 of Ortega (1990) asserts that convergence of Algorithm 3 is guaranteed when \mathbf{x}_{init} is sufficiently close to \mathbf{x}^* , the components of \mathbf{g} are differentiable at \mathbf{x}^* and

$$\rho(\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x}^*)) < 1.$$

Here $\rho(\mathbf{A})$ denotes the *spectral radius* of the square matrix \mathbf{A} , defined to be

$$\rho(\mathbf{A}) \equiv \text{maximum of the absolute values of the eigenvalues of } \mathbf{A}.$$

Theorem 8.4 of Ackleh et al. (2010) provides a similar condition in terms of the *spectral norm* $\|\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x}^*)\|_{\text{spec}}$ where

$$\|\mathbf{A}\|_{\text{spec}} \equiv \sqrt{\text{largest eigenvalue of } \mathbf{A}^T \mathbf{A}}.$$

- There are also theorems that guarantee convergence of Algorithm 3 for particular choices of \mathbf{x}_{init} . If D_0 is a closed convex subset of D such that $\mathbf{g} : D_0 \rightarrow D_0$, the entries of $\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x})$ are each bounded and continuous on D_0 and

$$\sup_{\mathbf{x} \in D_0} \|\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x})\|_{\text{spec}} \leq \alpha < 1$$

then Algorithm 3 will converge from any initial point $\mathbf{x}_{\text{init}} \in D_0$ (Theorems 8.2 and 8.3 of Ackleh *et al.*, 2010).

Despite this elegant theory, it is difficult to apply in practice with regards to choosing \mathbf{g} and \mathbf{x}_{init} so that Algorithm 3 converges. This is exemplified in Section 4.2 when we return to the Example 1 optimization problem. We also note that $\|\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x})\|_{\text{spec}} < 1$ near \mathbf{x}^* is a *sufficient* but not *necessary* condition for convergence of fixed point iteration. Nevertheless, the function

$$\rho(\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x}))$$

is a useful convergence diagnostic for fixed-point iteration. For instance, if $\rho(\mathbf{D}_{\mathbf{x}} \mathbf{g}(\mathbf{x})) \gg 1$ during the iterations then this would indicate convergence problems and the possibility of non-existence of a fixed point \mathbf{x}^* .

Various adjustments to fixed-point iteration have been proposed to enhance convergence. For example, in the context of semiparametric mean field variational Bayes, Section 7 of Minka & Knowles (2011) describes a *damping* adjustment.

3.1.1 NEWTON-RAPHSON ITERATION

Newton-Raphson iteration is a special case of fixed-point iteration for optimizing \mathbf{f} with the \mathbf{g} function taking the form

$$\mathbf{g}_{\text{NR}}(\mathbf{x}) = \mathbf{x} - \{\mathbf{H}_{\mathbf{x}} \mathbf{f}(\mathbf{x})\}^{-1} \mathbf{D}_{\mathbf{x}} \mathbf{f}(\mathbf{x})^T \quad (33)$$

where $\mathbf{H}_{\mathbf{x}} \mathbf{f}(\mathbf{x})$ denotes the Hessian matrix of $\mathbf{f}(\mathbf{x})$ as formally defined in Appendix A.1. Assuming existence of $\{\mathbf{H}_{\mathbf{x}} \mathbf{f}(\mathbf{x})\}^{-1}$, it is easily shown that $\mathbf{x} = \mathbf{g}_{\text{NR}}(\mathbf{x})$ if and only if $\mathbf{D}_{\mathbf{x}} \mathbf{f}(\mathbf{x})^T = \mathbf{0}$. This leads to Algorithm 4, which conveys the generic form of Newton-Raphson iteration.

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

Cycle:

$$\mathbf{x} \leftarrow \mathbf{x} - \{\mathbf{H}_{\mathbf{x}} \mathbf{f}(\mathbf{x})\}^{-1} \mathbf{D}_{\mathbf{x}} \mathbf{f}(\mathbf{x})^T$$

until convergence.

Algorithm 4: *The Newton-Raphson algorithm in generic form.*

Some pertinent features of Algorithm 4 are:

- The function \mathbf{g}_{NR} in (33) has the property

$$\rho(\mathbf{D}_{\mathbf{x}} \mathbf{g}_{\text{NR}}(\mathbf{x}^*)) = 0 \quad (34)$$

for stationary points \mathbf{x}^* . A proof is given in Appendix A.4. Therefore, via Theorem 8.4 of Ackleh et al. (2010), convergence to \mathbf{x}^* is guaranteed from a sufficiently close \mathbf{x}_{init} .

- If \mathbf{x}_{init} is such that Algorithm 4 is convergent to \mathbf{x}^* then, under certain regularity conditions, convergence is *quadratic*, in that the number of significant figures doubles on each iteration.

- Locating \mathbf{x}_{init} values sufficiently close to \mathbf{x}^* for convergence to occur can be difficult in practice and it is common to combine Newton-Raphson iteration with more robust optimization strategies, such as the Nelder-Mead simplex method.
- A disadvantage of Newton-Raphson iteration compared with other fixed-point iterative schemes is the requirement for second order partial derivatives, corresponding to the entries of the Hessian matrix $\mathbf{H}_{\mathbf{x}} f(\mathbf{x})$. A feeling for the type of additional calculus needed is given in Appendix A.7.
- A variant of Newton-Raphson optimization known as *damped* Newton-Raphson employs line searches (or backtracking) in order to achieve much improved convergence behavior. See, e.g., Section 9.5.2 of Boyd and Vandenberghe (2004).

3.2 Nonlinear Conjugate Gradient Method

The *nonlinear conjugate gradient method* is based on the *conjugate gradient method*, an established iterative approach to solving large systems of linear equations (Hestenes and Stiefel, 1952). The former arises from applying the latter to the linear system that arises from optimization of a multivariate quadratic function. Details of the nonlinear conjugate gradient method are given in Section 10.8 of Press et al. (2007). Algorithm 5 lists the *Polak-Ribière* version of the nonlinear conjugate gradient method for *maximization* of f over D . Since minimization of f is equivalent to maximization of $-f$ it is straightforward to adapt Algorithm 5 to the minimization problem. We choose the Polak-Ribière form here, but another one is the *Fletcher-Reeves* form given by $\beta \leftarrow (\mathbf{v}_{\text{curr}}^T \mathbf{v}_{\text{curr}})/(\mathbf{v}_{\text{prev}}^T \mathbf{v}_{\text{prev}})$.

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

$$\mathbf{v}_{\text{prev}} \leftarrow \mathbf{D}_{\mathbf{x}} f(\mathbf{x})^T \quad ; \quad \alpha \leftarrow \underset{\alpha > 0}{\operatorname{argmax}} f(\mathbf{x} + \alpha \mathbf{v}_{\text{prev}})$$

$$\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{v}_{\text{prev}} \quad ; \quad \mathbf{s} \leftarrow \mathbf{v}_{\text{prev}}$$

Cycle:

$$\mathbf{v}_{\text{curr}} \leftarrow \mathbf{D}_{\mathbf{x}} f(\mathbf{x})^T \quad ; \quad \beta \leftarrow \mathbf{v}_{\text{curr}}^T (\mathbf{v}_{\text{curr}} - \mathbf{v}_{\text{prev}}) / (\mathbf{v}_{\text{prev}}^T \mathbf{v}_{\text{prev}})$$

$$\mathbf{s} \leftarrow \beta \mathbf{s} + \mathbf{v}_{\text{curr}} \quad ; \quad \alpha \leftarrow \underset{\alpha > 0}{\operatorname{argmax}} f(\mathbf{x} + \alpha \mathbf{s})$$

$$\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{s} \quad ; \quad \mathbf{v}_{\text{prev}} \leftarrow \mathbf{v}_{\text{curr}}$$

until convergence.

Algorithm 5: *The nonlinear conjugate gradient method for maximization of the function f with the Polak-Ribière form of the β parameter.*

A key aspect of the nonlinear conjugate gradient method is that, on each iteration, it takes a step in the direction $\mathbf{D} f(\mathbf{x})^T$ from the current position at \mathbf{x} . This is the steepest instantaneous direction on the f surface. The parameter denoted by β has several alter-

native forms. Nonlinear conjugate gradient methods have been shown to have good global convergence properties (Dai and Yuan, 1999).

3.3 Other Optimization Strategies

Other popular optimization strategies include *ascent (or descent) algorithms* (e.g. Boyd and Vandenberghe, 2004, Section 9.3), *quasi-Newton methods* (e.g. Givens and Hoeting, 2005, Section 2.2.2.3), the *Gauss-Newton method* (e.g. Givens and Hoeting, 2005, Section 2.2.3), *stochastic gradient descent* (e.g. Bottou, 2004) and the *Nelder-Mead simplex method* (Nelder and Mead, 1965). The last of these has the attraction of not requiring derivatives of f and is generally more robust than derivative-based methods.

3.4 Application to Semiparametric Mean Field Variational Bayes

We now focus on the optimization component of Algorithm 2

$$\operatorname{argmax}_{\xi' \in \Xi} \left\{ \log p(\mathcal{D}; q, \xi')^{[\phi]} \right\} \quad (35)$$

and discuss ways in which numerical optimization strategies described in Sections 3.1–3.3 are applicable.

The stationary condition for the maximizer in (35) is

$$D_{\xi} \log p(\mathcal{D}; q, \xi)^{[\phi]T} = \mathbf{0}$$

and this may be manipulated in any of a number of ways to produce an equation of the form $\xi = \mathbf{g}(\xi)$ for some function \mathbf{g} . Fixed-point iteration Algorithm 3 can then be entertained but, as discussed in Section 3.1, converge is not guaranteed for arbitrary \mathbf{g} . We study this issue in the context of Examples 1 and 2 in Sections 4.2 and 4.3.

Newton-Raphson iteration involves iterative updates:

$$\xi \leftarrow \xi - \{H_{\xi} \log p(\mathcal{D}; q, \xi)^{[\phi]}\}^{-1} D_{\xi} \log p(\mathcal{D}; q, \xi)^{[\phi]T}$$

and so is a candidate for insertion into Algorithm 2 for updating the pre-specified parametric q -density parameters.

Another alternative is, of course, updating ξ according to one or more iterations of the nonlinear conjugate gradient method given by Algorithm 5, or any other iterative optimization scheme. However, convergence needs to be monitored. For high-dimensional ξ , *speed* of convergence may be also be an important factor. Next we discuss an adjustment aimed at improving the convergence speed of gradient-based algorithms.

3.4.1 RIEMANNIAN GEOMETRY ADJUSTMENT

As explained in, for example, Section 6.2 of Murray and Rice (1993) the density function family $\{q(\phi; \xi) : \xi \in \Xi\}$ can be viewed as a *submanifold* of a *Riemannian manifold*. Riemannian manifolds do not necessarily have a *flat* Euclidean geometry. For example, the Riemannian manifold corresponding to the Univariate Normal family:

$$\left\{ \frac{1}{\sqrt{2\pi\sigma_{q(\phi)}^2}} \exp \left\{ -\frac{(\phi - \mu_{q(\phi)})^2}{2\sigma_{q(\phi)}^2} \right\} : \mu_{q(\phi)} \in \mathbb{R}, \sigma_{q(\phi)}^2 > 0 \right\} \quad (36)$$

has *hyperbolic geometry* (Murray and Rice, 1993, Example 6.6.2) which is *curved*. Therefore notions such as closeness of two members of (36) and steepness of gradients when searching over the parameter space $\Xi = \mathbb{R} \times \mathbb{R}_+$ are not properly captured by the Euclidean geometry notions of distance and slope. Adjustments for the Riemannian geometry of the family often improve convergence of optimization algorithms for solving problems such as (35). More detailed discussion on this issue is given in Section 2.2 of Honkela et al. (2010) and Section 2.3 of Hoffman et al. (2013).

Consider an optimization method that uses gradients of the form

$$D_{\xi} \log \underline{p}(\mathcal{D}; q, \xi)^{[\phi]T}$$

to find the direction of steepest descent of the objective function $\log \underline{p}(\mathcal{D}; q, \xi)^{[\phi]}$. The Riemannian geometry adjustment is to instead use

$$[-E\{H_{\xi} \log q(\phi; \xi)\}]^{-1} D_{\xi} \log \underline{p}(\mathcal{D}; q, \xi)^{[\phi]T} \quad (37)$$

where the premultiplying matrix is the inverse *Fisher information* of $q(\phi; \xi)$. In the Machine Learning literature (e.g. Amari, 1998) (37) is often labeled the *natural* or *Riemannian gradient* of $\log \underline{p}(\mathcal{D}; q, \xi)^{[\phi]}$ with respect to ξ and the corresponding geometry is called *information geometry*. If $q(\phi; \xi)$ corresponds to the Univariate Normal family (36) then the Fisher information matrix is $\text{diag}(\sigma_{q(\phi)}^{-2}, \frac{1}{2}\sigma_{q(\phi)}^{-4})$. Therefore, from (37), the natural gradient of $\log \underline{p}(q; \mu_{q(\phi)}, \sigma_{q(\phi)}^2)^{[\phi]}$ with respect to $(\mu_{q(\phi)}, \sigma_{q(\phi)}^2)$ is given by

$$\begin{bmatrix} \sigma_{q(\phi)}^2 \frac{\partial \underline{p}(q; \mu_{q(\phi)}, \sigma_{q(\phi)}^2)^{[\phi]}}{\partial \mu_{q(\phi)}} \\ 2\sigma_{q(\phi)}^4 \frac{\partial \underline{p}(q; \mu_{q(\phi)}, \sigma_{q(\phi)}^2)^{[\phi]}}{\partial \sigma_{q(\phi)}^2} \end{bmatrix}^T. \quad (38)$$

Honkela et al. (2010) is a major contribution to semiparametric mean field variational Bayes methodology and their Algorithm 1 uses the nonlinear conjugate gradient method (Algorithm 5) with natural gradients rather than ordinary gradients. Via both simple examples and numerical studies, they make a compelling case for the use of natural gradients for optimization of the parameters of the pre-specified parametric q -density function.

3.5 Summary of Semiparametric Mean Field Variational Bayes Ramifications

In this section we have discussed several iterative numerical optimization strategies. Any of these are candidates for the updating ξ in the Algorithm 2 cycle. Special mention has been given to the well-known Newton-Raphson iteration since it can achieve very rapid convergence and the nonlinear conjugate gradient method which has been shown to be effective in semiparametric mean field variational Bayes contexts when the Riemannian geometry adjustment of Section 3.4.1 is employed (Honkela et al., 2010).

General fixed-point iteration has been discussed in detail. It has the advantage of yielding particularly simple iterative updates for ξ in Algorithm 2. Established theory shows that the spectral radius of the derivative matrix of the fixed-point function can be

used to assess the quality of the scheme. In Section 4 we will explain how a particular fixed-point iteration scheme, which we call *natural* fixed-point iteration, has attractive properties when $q(\cdot; \boldsymbol{\xi})$ is an exponential family density function. We will also revisit Examples 1 and 2 in Section 4 and make some comparisons among various numerical optimization strategies. Natural fixed-point iteration is seen to perform particularly well.

4. Exponential Family Special Case

We now focus on the important special case where the parametric density function family $\{q(\boldsymbol{\phi}; \boldsymbol{\xi}) : \boldsymbol{\xi} \in \Xi\}$ can be expressed in exponential family form:

$$q(\boldsymbol{\phi}; \boldsymbol{\eta}) = \exp\{\mathbf{T}(\boldsymbol{\phi})^T \boldsymbol{\eta} - A(\boldsymbol{\eta})\} h(\boldsymbol{\phi}), \quad \boldsymbol{\eta} \in H, \quad (39)$$

where $\boldsymbol{\eta}$ is a one-to-one transformation of $\boldsymbol{\xi}$ and H is the image of Ξ under this transformation. In (39) $A(\boldsymbol{\eta})$ is called the *log-partition function* and $h(\boldsymbol{\phi})$ is called the *base measure*. For example, the Univariate Normal density function family used in Example 1:

$$q(\phi; \mu_{q(\phi)}, \sigma_{q(\phi)}^2) = \frac{1}{\sqrt{2\pi\sigma_{q(\phi)}^2}} \exp\left\{-\frac{(\phi - \mu_{q(\phi)})^2}{2\sigma_{q(\phi)}^2}\right\}, \quad \mu_{q(\phi)} \in \mathbb{R}, \sigma_{q(\phi)}^2 > 0,$$

can be expressed as (39) with

$$\mathbf{T}(\phi) = \begin{bmatrix} \phi \\ \phi^2 \end{bmatrix}, \quad \boldsymbol{\eta} \equiv \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} \mu_{q(\phi)}/\sigma_{q(\phi)}^2 \\ -1/(2\sigma_{q(\phi)}^2) \end{bmatrix}, \quad A(\boldsymbol{\eta}) = -\frac{1}{4} \eta_1^2/\eta_2 - \frac{1}{2} \log(-2\eta_2)$$

and $h(\phi) = (2\pi)^{-1/2}$. The natural parameter space is $H = \{(\eta_1, \eta_2) : \eta_1 \in \mathbb{R}, \eta_2 < 0\}$. Even though semiparametric mean field variational Bayes can involve pre-specification of an arbitrary parametric family, virtually all methodology and examples in the existing literature involves pre-specification within an exponential family. Exponential family distributions also play an important role in the theory of mean field variational Bayes (e.g. Sato, 2001; Beal and Ghahramani, 2006; Wainwright and Jordan, 2008).

Now consider the general factor graph set-up described in Section 2.2.1 with the approximate marginal log-likelihood $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\eta})^{[\phi]}$ given by (13) but as a function of the natural parameter vector $\boldsymbol{\eta}$. Define

$$\text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\} \equiv \sum_{j \in \text{neighbors}(\phi)} E_{q(\boldsymbol{\phi}; \boldsymbol{\eta})}\{\log(f_j)\}$$

so that

$$\log \underline{p}(\mathcal{D}; q, \boldsymbol{\eta})^{[\phi]} = \text{Entropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\} + \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}.$$

An advantage of working with exponential family density functions is that the entropy takes the simple form:

$$\text{Entropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\} = A(\boldsymbol{\eta}) - \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) \boldsymbol{\eta} + E[\exp\{h(\boldsymbol{\phi})\}]$$

where the notation of Appendix A.2 is being used. Moreover, as shown in Lemma 1 of Appendix A.5, the derivative vector of $\text{Entropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}$ is simply

$$\mathbf{D}_{\boldsymbol{\eta}} \text{Entropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\} = -\boldsymbol{\eta}^T \mathbf{H}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}). \quad (40)$$

This implies that the stationary point condition

$$\{D_{\boldsymbol{\eta}} \log \underline{p}(q; \boldsymbol{\eta})^{[\phi]}\}^T = \mathbf{0} \quad (41)$$

is equivalent to

$$\boldsymbol{\eta} = \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} D_{\boldsymbol{\eta}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}^T. \quad (42)$$

Algorithm 1 of Knowles and Minka (2011) is a fixed-point iteration scheme based on (42).

One further interesting and useful connection concerns the *mean parameter* vector

$$\boldsymbol{\tau} \equiv E\{T(\boldsymbol{\phi})\}$$

which is related to the natural parameter vector via

$$\boldsymbol{\tau} = D_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T.$$

Under suitable technical conditions $\boldsymbol{\tau}$ is a one-to-one transformation of $\boldsymbol{\eta}$. Also the chain rule for differentiation of a smooth function s , listed as Lemma 2 in Appendix A.5, is

$$D_{\boldsymbol{\eta}} s = (D_{\boldsymbol{\tau}} s)(D_{\boldsymbol{\eta}} \boldsymbol{\tau}) = (D_{\boldsymbol{\tau}} s) D_{\boldsymbol{\eta}} \{D_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T\} = (D_{\boldsymbol{\tau}} s) H_{\boldsymbol{\eta}} A(\boldsymbol{\eta}).$$

which leads to

$$D_{\boldsymbol{\tau}} s = \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} D_{\boldsymbol{\eta}} s. \quad (43)$$

Putting all of these relationships together we have:

Result 1 *Let $\boldsymbol{\xi}$ be an arbitrary differentiable one-to-one transformation of $\boldsymbol{\eta}$. Then the stationary point condition (41) is equivalent to each the following conditions:*

- (a) $\boldsymbol{\eta} = \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} D_{\boldsymbol{\eta}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}^T$,
- (b) $\boldsymbol{\eta} = \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} (D_{\boldsymbol{\eta}} \boldsymbol{\xi})^T D_{\boldsymbol{\xi}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\xi})\}^T$,
- (c) $\boldsymbol{\eta} = D_{\boldsymbol{\tau}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\tau})\}^T$ and
- (d) $\boldsymbol{\eta} = (D_{\boldsymbol{\tau}} \boldsymbol{\xi})^T D_{\boldsymbol{\xi}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\xi})\}^T$.

We make the following remarks concerning Result 1:

- Result 1(a) immediately gives rise to the following fixed-point iteration scheme in the natural parameter space $\boldsymbol{\eta} \in H$:

$$\boldsymbol{\eta} \leftarrow \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} D_{\boldsymbol{\eta}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}^T. \quad (44)$$

We refer to (44) as the *natural fixed-point iteration scheme* and denote the corresponding fixed-point function by

$$\mathbf{g}_{\text{nat}}(\boldsymbol{\eta}) \equiv \{H_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\}^{-1} D_{\boldsymbol{\eta}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}^T.$$

According to the theory of fixed-point iteration discussed in Section 3.1, convergence of (44) is implied by

$$\rho(D_{\boldsymbol{\eta}} \mathbf{g}_{\text{nat}}(\boldsymbol{\eta})) < 1$$

in a neighborhood of the maximizer $\boldsymbol{\eta}^*$.

- Result 1(b)–(d) offer the possibility of more convenient forms for the fixed-point updates in terms of derivatives of the common parameters or mean parameters. Particularly noteworthy is the fact that Result 1(c)–(d) do not require computation of $\{\mathbf{H}_\eta A(\boldsymbol{\eta})\}^{-1}$. We make use of this situation for the Multivariate Normal family in Section 4.1.
- The Fisher information of $q(\boldsymbol{\phi} : \boldsymbol{\eta})$ is

$$-E\{\mathbf{H}_\eta \log q(\boldsymbol{\phi}; \boldsymbol{\eta})\} = \mathbf{H}_\eta A(\boldsymbol{\eta})$$

which implies that the natural fixed-point iteration scheme (44) involves updating $\boldsymbol{\eta}$ according to natural Riemannian gradients of $\text{NonEntropy}(q; \boldsymbol{\tau})$. From Result 1(c), an equivalent updating scheme is

$$\boldsymbol{\eta} \leftarrow \mathbf{D}_\tau \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\tau})\}^T$$

which simply involves updating $\boldsymbol{\eta}$ according to the direction of maximum slope on the $\text{NonEntropy}(q; \boldsymbol{\tau})$ surface in the $\boldsymbol{\tau}$ space.

- The forms for the stationary point in Result 1 can also be used to derive iterative Newton-Raphson schemes for maximizing $\log p(q; \boldsymbol{\eta})^{[\phi]}$. An example, corresponding to Result 1(a) and optimization within the $\boldsymbol{\eta}$ space, is

$$\begin{aligned} \boldsymbol{\eta} \leftarrow \boldsymbol{\eta} - & [\mathbf{H}_\eta \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\} - \mathbf{H}_\eta A(\boldsymbol{\eta}) - (\boldsymbol{\eta}^T \otimes \mathbf{I}) \mathbf{D}_\eta \text{vec}\{\mathbf{H}_\eta A(\boldsymbol{\eta})\}]^{-1} \\ & \times [\mathbf{D}_\eta \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}^T - \mathbf{H}_\eta A(\boldsymbol{\eta}) \boldsymbol{\eta}]. \end{aligned}$$

The vec operator flattens a square matrix into a column vector and is defined formally in Appendix A.1.

Any of the other optimization methods mentioned in Section 3 can also be applied to the problem of obtaining

$$\boldsymbol{\eta}^* \equiv \underset{\boldsymbol{\eta} \in H}{\text{argmax}} \{\log p(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]}\} = \underset{\boldsymbol{\eta} \in H}{\text{argmax}} [A(\boldsymbol{\eta}) - \mathbf{D}_\eta A(\boldsymbol{\eta}) \boldsymbol{\eta} + \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\eta})\}]$$

and those involving gradients benefit from the entropy derivative result (40). Additionally, relationship (43) implies that natural (Riemannian) gradients of the objective function in the natural parameter space are equivalent to ordinary Euclidean gradients in the mean parameter space.

4.1 Multivariate Normal Special Case

We now focus on the important special case of $q(\boldsymbol{\phi}; \boldsymbol{\xi})$ being a d -variate Multivariate Normal density function:

$$q(\boldsymbol{\phi}; \boldsymbol{\mu}_{q(\boldsymbol{\phi})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\phi})}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}_{q(\boldsymbol{\phi})}|^{-1/2} \exp\{-\frac{1}{2}(\boldsymbol{\phi} - \boldsymbol{\mu}_{q(\boldsymbol{\phi})})^T \boldsymbol{\Sigma}_{q(\boldsymbol{\phi})}^{-1} (\boldsymbol{\phi} - \boldsymbol{\mu}_{q(\boldsymbol{\phi})})\}.$$

Let

$$\boldsymbol{\xi} \equiv \begin{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\phi})} \\ \text{vec}(\boldsymbol{\Sigma}_{q(\boldsymbol{\phi})}) \end{bmatrix}$$

be the vector of common parameters. An explicit form for the natural fixed point iteration updates in terms of $\boldsymbol{\mu}_{q(\phi)}$ and $\boldsymbol{\Sigma}_{q(\phi)}$ was derived by Wand (2014) and appears as equation (7) there. However Result 1 affords a more direct derivation of the same result, that benefits from (43) and the cancellation of the $\mathbf{H}_\eta A(\boldsymbol{\eta})$ matrix. We can also obtain a neater alternative explicit form by using a differentiation identity, given as Lemma 4 in Appendix A.5. The essence of Lemma 4 is given in the appendix of Opper and Archambeau (2009).

Result 2 *Natural fixed-point iteration for $q(\phi; \boldsymbol{\xi})$ corresponding to the $N(\boldsymbol{\mu}_{q(\phi)}, \boldsymbol{\Sigma}_{q(\phi)})$ density function is equivalent to the following updating scheme:*

$$\begin{cases} \mathbf{v}_{q(\phi)} \leftarrow \mathbf{D}_{\boldsymbol{\mu}_{q(\phi)}} \text{NonEntropy}(q; \boldsymbol{\mu}_{q(\phi)}, \boldsymbol{\Sigma}_{q(\phi)})^T \\ \boldsymbol{\Sigma}_{q(\phi)} \leftarrow -\{\mathbf{H}_{\boldsymbol{\mu}_{q(\phi)}} \text{NonEntropy}(q; \boldsymbol{\mu}_{q(\phi)}, \boldsymbol{\Sigma}_{q(\phi)})\}^{-1} \\ \boldsymbol{\mu}_{q(\phi)} \leftarrow \boldsymbol{\mu}_{q(\phi)} + \boldsymbol{\Sigma}_{q(\phi)} \mathbf{v}_{q(\phi)}. \end{cases}$$

Appendix A.6 provides details on how Result 2 follows from Result 1.

Result 2 facilitates a semiparametric mean field variational Bayes algorithm that requires only the first and second order derivatives of $\text{NonEntropy}(q; \boldsymbol{\mu}_{q(\phi)}, \boldsymbol{\Sigma}_{q(\phi)})$ with respect to $\boldsymbol{\mu}_{q(\phi)}$. Concrete illustrations are given in Section 4.3 and Appendix A of Menictas and Wand (2015).

In the case where $q(\phi; \boldsymbol{\xi})$ is the Univariate Normal density function with mean $\mu_{q(\phi)}$ and variance $\sigma_{q(\phi)}^2$ Result 2 leads to the following common parameter updates for the natural fixed point iterative scheme:

$$\begin{cases} v_{q(\phi)} \leftarrow \frac{\partial \text{NonEntropy}(q; \mu_{q(\phi)}, \sigma_{q(\phi)}^2)}{\partial \mu_{q(\phi)}} \\ \sigma_{q(\phi)}^2 \leftarrow -1 / \left\{ \frac{\partial^2 \text{NonEntropy}(q; \mu_{q(\phi)}, \sigma_{q(\phi)}^2)}{\partial \mu_{q(\phi)}^2} \right\} \\ \mu_{q(\phi)} \leftarrow \mu_{q(\phi)} + \sigma_{q(\phi)}^2 v_{q(\phi)}. \end{cases} \quad (45)$$

Despite its use of natural gradients, there is no automatic guarantee that iteration of (45) leads to convergence to the maximum of $\log \underline{p}(\mathcal{D}; q, \boldsymbol{\xi})^{[\phi]}$ on any given cycle of Algorithm 2. However, the fixed-point iteration theory summarized in Section 3.1 provides some guidance. We now use Example 1 to illustrate this point using the natural fixed-point iteration scheme (45), an alternative simpler fixed-point scheme and a Newton-Raphson scheme.

4.2 Application to Example 1

Consider again the Gumbel random sample example introduced in Section 2.2.2 and the problem of minimum Kullback-Leibler approximation of $p(\phi|\mathbf{x})$ within the Univariate Normal family. As shown there, the optimization problem is encapsulated in (19) and (20).

The Newton-Raphson scheme that arises directly from (19) is

$$\begin{bmatrix} \mu_{q(\phi)} \\ \sigma_{q(\phi)}^2 \end{bmatrix} \leftarrow \begin{bmatrix} \mu_{q(\phi)} \\ \sigma_{q(\phi)}^2 \end{bmatrix} - \{\mathbf{H}_{\text{Ex1}}^{fN}(\mu_{q(\phi)}, \sigma_{q(\phi)}^2; n, \sum_{i=1}^n e^{-x_i}, \mu_\phi, \sigma_\phi^2)\}^{-1} \times \mathbf{D}_{\text{Ex1}}^{fN}(\mu_{q(\phi)}, \sigma_{q(\phi)}^2; n, \sum_{i=1}^n e^{-x_i}, \mu_\phi, \sigma_\phi^2)^T. \quad (46)$$

Differentiation with respect to $(\mu_\phi, \sigma_\phi^2)$ is suppressed in the \mathbf{D} and \mathbf{H} on the right-hand side of (46). Simple calculus leads to (46) being equivalent to the fixed-point iterative scheme

$$\begin{bmatrix} \mu_{q(\phi)} \\ \sigma_{q(\phi)}^2 \end{bmatrix} \leftarrow \mathbf{g}_{\text{NR}} \left(\begin{bmatrix} \mu_{q(\phi)} \\ \sigma_{q(\phi)}^2 \end{bmatrix}; n, \sum_{i=1}^n e^{-x_i}, \mu_\phi, \sigma_\phi^2 \right) \quad (47)$$

where

$$\begin{aligned} \mathbf{g}_{\text{NR}} \left(\begin{bmatrix} x \\ y \end{bmatrix}; a, b, c, d \right) &\equiv \begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} -b e^{x+\frac{1}{2}y} - 1/d & -\frac{1}{2} b e^{x+\frac{1}{2}y} \\ -\frac{1}{2} b e^{x+\frac{1}{2}y} & -\frac{1}{2y^2} - \frac{1}{4} b e^{x+\frac{1}{2}y} \end{bmatrix}^{-1} \\ &\quad \times \begin{bmatrix} a - b e^{x+\frac{1}{2}y} - (x-c)/d \\ \frac{1}{2y} - \frac{1}{2} b e^{x+\frac{1}{2}y} - 1/(2d) \end{bmatrix}. \end{aligned}$$

According to (45), the natural fixed-point iteration scheme (44) takes the form (47), but with \mathbf{g}_{NR} replaced by \mathbf{g}_{nat} where

$$\mathbf{g}_{\text{nat}} \left(\begin{bmatrix} x \\ y \end{bmatrix}; a, b, c, d \right) \equiv \begin{bmatrix} x + \{a - b e^{x+\frac{1}{2}y} - (x-c)/d\} / (b e^{x+\frac{1}{2}y} + d^{-1}) \\ 1 / (b e^{x+\frac{1}{2}y} + d^{-1}) \end{bmatrix}.$$

Lastly, there is the very simple fixed-point iteration scheme that arises from full simplification of $\mathbf{D} f_{\text{Ex1}}^N(\mu_{q(\phi)}, \sigma_{q(\phi)}^2; n, \sum_{i=1}^n e^{-x_i}, \mu_\phi, \sigma_\phi^2)^T = \mathbf{0}$, and corresponds to fixed points of

$$\mathbf{g}_{\text{simp}} \left(\begin{bmatrix} x \\ y \end{bmatrix}; a, b, c, d \right) \equiv \begin{bmatrix} c + d(a - b e^{x+\frac{1}{2}y}) \\ 1 / (b e^{x+\frac{1}{2}y} + d^{-1}) \end{bmatrix}.$$

In Figure 4 we compare \mathbf{g}_{NR} , \mathbf{g}_{nat} and \mathbf{g}_{simp} in terms of the behavior of the spectral norm function $\rho(\mathbf{D}\mathbf{g}(x, y))$ and convergence of the fixed-point iterative scheme. We simulated data from the $n = 20$ version of the Gumbel random sample model (14) with the value of ϕ set to 0. The sufficient statistic $\sum_{i=1}^{20} \exp(-x_i)$ fully determines f_{Ex1}^N and has a mean of 20. In an effort to exhibit typical behavior, we selected a sample that produced a sufficient statistic value close to this mean. The actual value is $\sum_{i=1}^{20} \exp(-x_i) \approx 19.94$. The hyperparameters were set to $\mu_\phi = 0$ and $\sigma_\phi^2 = 10^{10}$. The optimal parameters in the minimum Kullback-Leibler Univariate Normal approximation to $p(\phi|\mathbf{x})$ are $(\mu_{q(\phi)}^*, (\sigma_{q(\phi)}^2)^*) \approx (0.2260, 0.0500)$. We set up a 101×101 pixel mesh of $(\mu_{q(\phi)}, \log(\sigma_{q(\phi)}^2))$ values around this optimum with limits $(\mu_{q(\phi)}^* - 5, \mu_{q(\phi)}^* + 5)$ and $(\log\{(\sigma_{q(\phi)}^*/5)^2\}, \log\{(5\sigma_{q(\phi)}^*)^2\})$. The upper panels of Figure 4 show the

$$\text{indicator of } \rho(\mathbf{D}\mathbf{g}(\mu_{q(\phi)}, \sigma_{q(\phi)}^2)) < 1 \text{ for } \mathbf{g} \in \{\mathbf{g}_{\text{NR}}, \mathbf{g}_{\text{nat}}, \mathbf{g}_{\text{simp}}\}.$$

The lower panels show the

$$\text{indicator of fixed-point iteration converging when starting from } (\mu_{q(\phi)}, \sigma_{q(\phi)}^2).$$

The top half of Figure 4 shows, via dark grey shading, that both $\rho(\mathbf{D}\mathbf{g}_{\text{NR}}(\mu_{q(\phi)}, \sigma_{q(\phi)}^2))$ and $\rho(\mathbf{D}\mathbf{g}_{\text{nat}}(\mu_{q(\phi)}, \sigma_{q(\phi)}^2))$ are below 1 in regions around the root. The dark grey region for

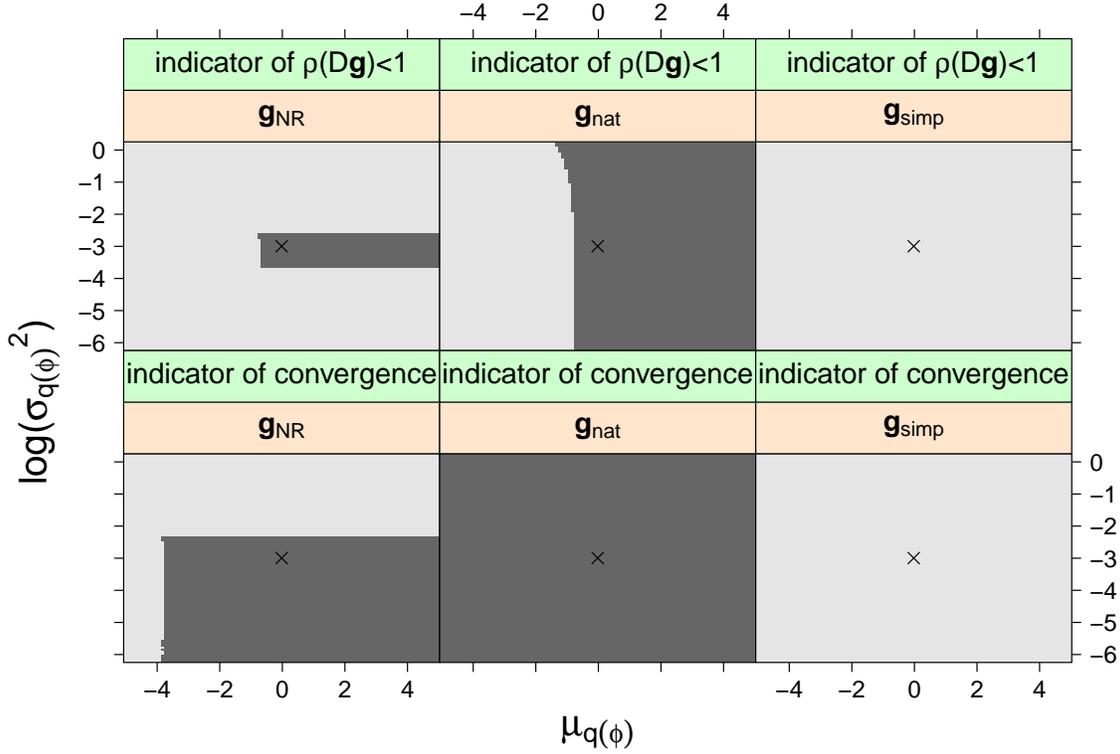


Figure 4: *Upper panels: Dark grey shading shows points where $\rho(\mathbf{Dg}(\mu_{\mathbf{q}(\phi)}, \sigma_{\mathbf{q}}^2)) < 1$ for $\mathbf{g} \in \{\mathbf{g}_{\text{NR}}, \mathbf{g}_{\text{nat}}, \mathbf{g}_{\text{simp}}\}$. Lower panels: Dark grey shading shows points from which fixed-point iteration, based on $\mathbf{g} \in \{\mathbf{g}_{\text{NR}}, \mathbf{g}_{\text{nat}}, \mathbf{g}_{\text{simp}}\}$, converges if initialized from that point. The optimum is shown by a cross in each panel and corresponds to minimum Kullback-Leibler divergence for a single $n = 20$ sample of the Gumbel random sample model with hyperparameters set to $\mu_{\phi} = 0$ and $\sigma_{\phi}^2 = 10^{10}$.*

\mathbf{g}_{nat} is much larger than that of \mathbf{g}_{NR} , suggesting that the former has better convergence properties according to the theory described in Section 3.1. The lower panels confirm this, with \mathbf{g}_{nat} -based fixed-point iteration converging from every initial value on the pixel mesh, but Newton-Raphson iteration not converging from the sub-region on the top and left-hand side of the mesh. Also note that $\rho(\mathbf{Dg}_{\text{simp}}(\mu_{\mathbf{q}(\phi)}, \sigma_{\mathbf{q}}^2)) \geq 1$ over the whole pixel mesh and \mathbf{g}_{simp} -based fixed-point iteration does not converge, regardless of initial point.

Figure 5 shows the iteration trajectories from four different starting values and four iterative schemes based on the same data and hyperparameter values as used in Figure 4. Also shown in each panel is an image plot of the surface being maximized, with a logarithmic scale used for $\sigma_{\mathbf{q}(\phi)}^2$. In addition to the fixed-point iteration schemes based on \mathbf{g}_{NR} and \mathbf{g}_{nat} we include the nonlinear conjugate gradient method given in Algorithm 5 and

the Riemannian geometry adjustment involving the natural gradients given by (38). In most cases the iterations led to convergence to $(\mu_{q(\phi)}^*, (\sigma_{q(\phi)}^2)^*)$ and the first three iterates are plotted. However, Newton-Raphson failed to converge from the starting values in each of the upper panels and the subsequent iterates are outside of the image plot boundaries.

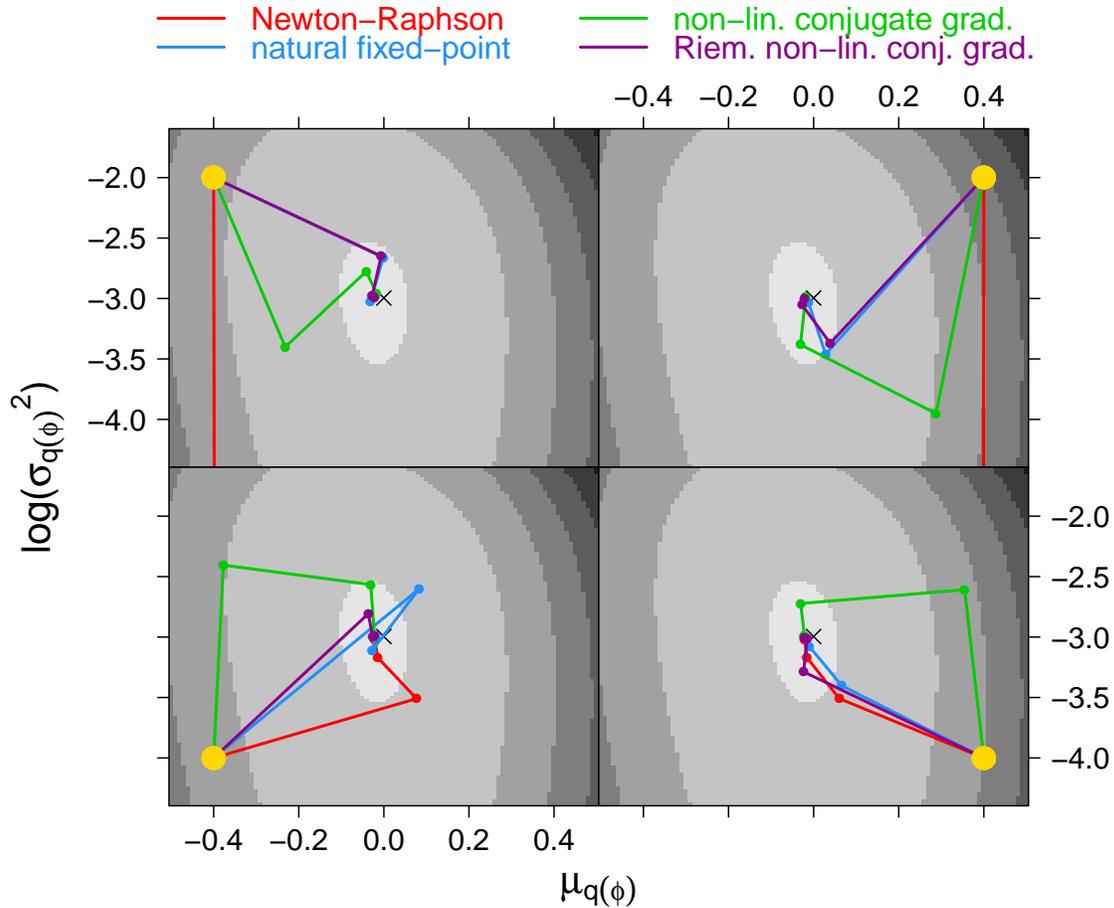


Figure 5: Iteration trajectories of four iterative algorithms aimed at solving the minimum Kullback-Leibler problem for a Gumbel random sample of size $n = 20$ with hyperparameters $\mu_{q(\phi)} = 0$ and $\sigma_{q(\phi)}^2 = 10^{10}$. The initial value differs for each panel and is shown by the yellow dot. The iterative algorithms are: (1) Newton-Raphson fixed-point iteration based on \mathbf{g}_{NR} , (2) natural fixed-point iteration based on \mathbf{g}_{nat} , (3) ordinary non-linear conjugate gradient method and (4) Riemannian non-linear conjugate gradient method.

The most striking feature of Figure 5 is the directness with which natural fixed-point iteration and the Riemannian non-linear conjugate gradient method converge from all four starting points and the similarity of their trajectories. This behavior is in keeping with the fact that both work with the more appropriate Riemannian gradients. The ordinary non-linear conjugate gradient trajectories are not as direct. Similar observations are made in Honkela et al. (2010). As demonstrated there, the payoffs from using Riemannian gradients in non-linear conjugate gradient updating are greater in higher-dimensional versions of semiparametric mean field variational Bayes. Based on Figure 5, we anticipate that natural fixed-point iteration is also very good in higher dimensions, and this is corroborated by experiments for Example 2 described in Section 4.3. Newton-Raphson fixed-point iteration is seen to be unreliable for this optimization problem and nowhere near as robust as natural fixed-point iteration. Lastly, we note that the behavior represented in Figures 4 and 5 persists across each of several other samples that we generated.

4.3 Application to Example 2

From (29) and some simple matrix algebra

$$\begin{aligned} & \text{NonEntropy}(q; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \\ &= \mathbf{y}^T \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} - \mathbf{1}^T \exp \left\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T) \right\} \\ & \quad - \frac{1}{2} \text{tr} \left(\begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \{ \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^T + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \} \right) \\ & \quad - \frac{1}{2} (p + K) \log(2\pi) - \frac{1}{2} p \log(\sigma_{\boldsymbol{\beta}}^2) - \frac{1}{2} K E_q \{ \log(\sigma^2) \} - \mathbf{1}^T \log(\mathbf{y}!) \end{aligned}$$

where

$$\mu_{q(1/\sigma^2)} = E_{q(1/\sigma^2)}(1/\sigma^2).$$

The derivatives appearing in Result 2 are

$$\begin{aligned} \mathbf{D}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \text{NonEntropy}(q; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})})^T &= \mathbf{C}^T \left[\mathbf{y} - \exp \{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T) \} \right] \\ & \quad - \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \end{aligned}$$

and

$$\begin{aligned} \mathbf{H}_{\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}} \text{NonEntropy}(q; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}) &= \\ & - \left(\mathbf{C}^T \text{diag} \{ \exp \{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T) \} \} \mathbf{C} + \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \right). \end{aligned}$$

It follows that the updates take the explicit form

$$\begin{cases} \mathbf{w}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \exp \{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T) \} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \left(\mathbf{C}^T \text{diag} \{ \mathbf{w}_{q(\boldsymbol{\beta}, \mathbf{u})} \} \mathbf{C} + \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \right)^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \left\{ \mathbf{C}^T (\mathbf{y} - \mathbf{w}_{q(\boldsymbol{\beta}, \mathbf{u})}) - \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \right\}. \end{cases}$$

This is equivalent to the fixed-point iteration scheme

$$\begin{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \\ \text{vech}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \end{bmatrix} \leftarrow \mathbf{g}_{\text{Ex2}} \left(\begin{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} \\ \text{vech}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}) \end{bmatrix}; \mathbf{y}, \mathbf{C}, \begin{bmatrix} \sigma_{\boldsymbol{\beta}}^{-2} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma^2)} \mathbf{I}_K \end{bmatrix} \right)$$

where

$$\mathbf{g}_{\text{Ex2}} \left(\begin{bmatrix} \boldsymbol{\mu} \\ \text{vech}(\boldsymbol{\Sigma}) \end{bmatrix}; \mathbf{y}, \mathbf{C}, \mathbf{M} \right) \equiv \begin{bmatrix} \boldsymbol{\mu} + [\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}]^{-1} \\ \quad \times [\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M} \boldsymbol{\mu}] \\ \text{vech} \left([\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}]^{-1} \right) \end{bmatrix}$$

and

$$\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \exp\{\mathbf{C} \boldsymbol{\mu} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma} \mathbf{C}^T)\}.$$

Note that the vech operator stores the unique entries of a symmetric matrix in a column vector. A formal definition of vech is given in Appendix A.1.

We simulated data according to the following special case of the Poisson mixed model:

$$\begin{aligned} y_{ij} | U_i &\sim \text{Poisson}\{\exp(\beta_0 + \beta_1 x_{ij} + U_i)\}, \quad U_i | \sigma^2 \sim N(0, \sigma^2), \\ 1 \leq i \leq m, \quad 1 \leq j \leq n, \quad \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}), \\ \sigma^2 | a &\sim \text{Inverse-Gamma}(\tfrac{1}{2}, 1/a), \quad a \sim \text{Inverse-Gamma}(\tfrac{1}{2}, A^{-2}). \end{aligned} \quad (48)$$

The hyperparameters were set at $\sigma_{\boldsymbol{\beta}} = A = 10^5$ and the sample sizes were $m = 30$, $n = 5$. Note that (48) is a special case of (21) with $\mathbf{Z} = \mathbf{I}_m \otimes \mathbf{1}_n$, where $\mathbf{1}_n$ is the $n \times 1$ vector with all entries equal to one. We then ran Algorithm 2 with $q(\boldsymbol{\beta}, \mathbf{u})$ pre-specified to be the $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ density function and a single natural fixed-point iteration in each cycle based on \mathbf{g}_{Ex2} . The fixed-point iteration search over values of $[\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}^T \text{vech}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})]^T$ is within an open subset of \mathbb{R}^{560} .

Figure 6 shows trace plots of $\log p(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ and $\rho(\mathbf{D} \mathbf{g}_{\text{Ex2}})$, based on the explicit expressions for $\mathbf{D} \mathbf{g}_{\text{Ex2}}$ given in Appendix A.7. The upper panel indicates that the algorithm becomes close to convergence after 6 – 10 iterations. After the same number of iterations the values of $\rho(\mathbf{D} \mathbf{g}_{\text{Ex2}})$ fall below 1 and settle at about 0.15.

Before leaving this example we note that numerical checks indicate that the optimal $\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}$ matrix is approximately sparse, with dominant diagonal entries. This implies the possibility of low-rank approximations to the above semiparametric mean field variational Bayes algorithm given, as described in Section 4.1.3 of Challis and Barber (2013).

5. A Non-Exponential Family Example

In the previous section it was seen that semiparametric mean field variational Bayes with $q(\boldsymbol{\phi}; \boldsymbol{\xi})$ having an exponential family form (39) leads to simplifications of the Kullback-Leibler minimization problem. However, $q(\boldsymbol{\phi}; \boldsymbol{\xi})$ does not have to be restricted in this way. In this section we illustrate semiparametric mean field variational Bayes with the q -density of $\boldsymbol{\phi}$ specified to be in a non-exponential family: the family of Skew-Normal density functions (Azzalini and Dalla Valle, 1996). Even though this family has a multivariate extension,

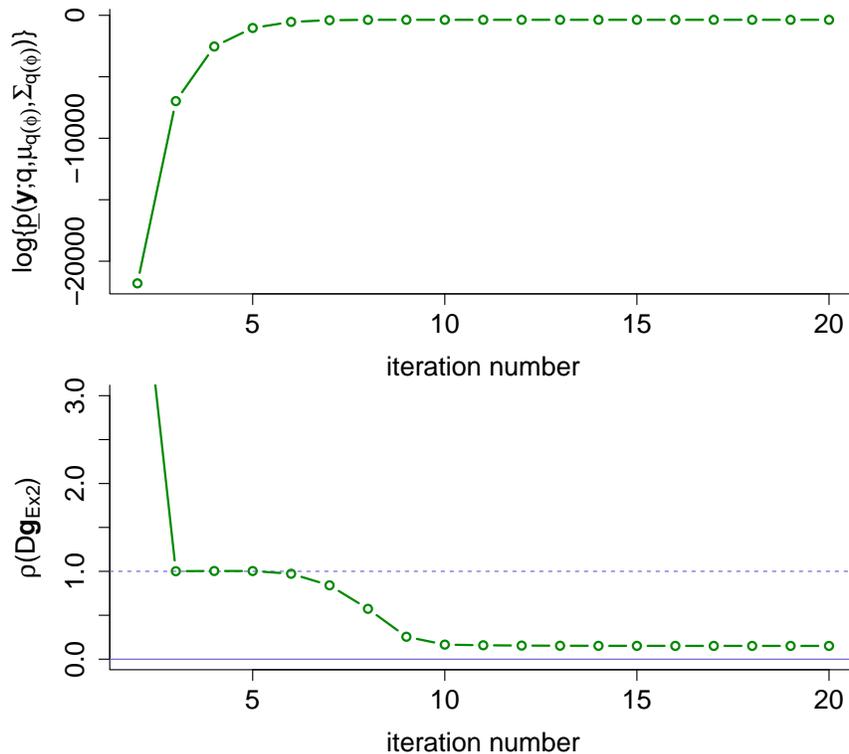


Figure 6: Trace plots of $\log p(\mathbf{y}; q, \boldsymbol{\xi})^{[(\beta, \mathbf{u})]}$ and $\rho(D\mathbf{g}_{\text{Ex2}})$ for the version of the Poisson mixed model given by (48) with sample sizes $m = 30$ and $n = 5$.

we restrict attention to the univariate case and, in particular, its use within the context of Example 1.

Specification of $q(\phi; \boldsymbol{\xi})$ being within the family of univariate Skew-Normal density functions entails having

$$q(\phi; \boldsymbol{\xi}) = \sqrt{\frac{2}{\pi\sigma_{q(\phi)}^2}} \exp\left\{-\frac{(\phi - \mu_{q(\phi)})^2}{2\sigma_{q(\phi)}^2}\right\} \Phi\left\{\frac{\lambda_{q(\phi)}(\phi - \mu_{q(\phi)})}{\sigma_{q(\phi)}}\right\} \quad (49)$$

where $\Phi(x) \equiv (2\pi)^{-1/2} \int_{-\infty}^x e^{-t^2/2} dt$ is the $N(0, 1)$ cumulative distribution function. The q -density parameter vector is $\boldsymbol{\xi} = (\mu_{q(\phi)}, \sigma_{q(\phi)}^2, \lambda_{q(\phi)})$ and the corresponding parameter space is $\boldsymbol{\Xi} = \mathbb{R} \times \mathbb{R}_+ \times \mathbb{R}$. Now consider the Example 1 setting with $q(\phi; \boldsymbol{\xi})$ restricted to the Skew-Normal family (49). The marginal log-likelihood lower bound, given by (17) and (18),

depends on the explicit expressions

$$E_{q(\phi;\boldsymbol{\xi})}(\phi) = \mu_{q(\phi)} + \frac{\sigma_{q(\phi)}\lambda_{q(\phi)}}{\sqrt{(\pi/2)(1+\lambda_{q(\phi)}^2)}}, \quad \text{Var}_{q(\phi;\boldsymbol{\xi})}(\phi) = \sigma_{q(\phi)}^2 \left\{ 1 - \frac{2\lambda_{q(\phi)}^2}{\pi(1+\lambda_{q(\phi)}^2)} \right\}$$

$$\text{and } M_{q(\phi;\boldsymbol{\xi})}(t) = 2 \exp\left(\mu_{q(\phi)} + \frac{1}{2}\sigma_{q(\phi)}^2 t^2\right) \Phi\left(\frac{\lambda_{q(\phi)}\sigma_{q(\phi)}t}{\sqrt{1+\lambda_{q(\phi)}^2}}\right)$$

where, as defined in Section 2.2.2, $M_{q(\phi;\boldsymbol{\xi})}$ is the moment generating function corresponding to $q(\phi;\boldsymbol{\xi})$. It also depends on

$$\text{Entropy}\{q(\phi;\boldsymbol{\xi})\} = \frac{1}{2}\{1 + \log(\pi/2) + \log(\sigma_{q(\phi)}^2)\} - \sqrt{\frac{2}{\pi}} \int_{-\infty}^{\infty} \log \Phi(\lambda_{q(\phi)}t) \Phi(\lambda_{q(\phi)}t) e^{-t^2/2} dt$$

which does not simplify any further. Plugging these expressions into (17) and (18) we get the Kullback-Leibler optimal Skew-Normal q -density function is $q(\phi; \mu_{q(\phi)}^*, (\sigma_{q(\phi)}^2)^*, \lambda_{q(\phi)}^*)$ where

$$\begin{bmatrix} \mu_{q(\phi)}^* \\ (\sigma_{q(\phi)}^2)^* \\ \lambda_{q(\phi)}^* \end{bmatrix} = \underset{\mu_{q(\phi)} \in \mathbb{R}, \sigma_{q(\phi)}^2 > 0, \lambda_{q(\phi)} \in \mathbb{R}}{\text{argmax}} \left\{ f_{\text{Ex1}}^{\text{SN}} \left(\mu_{q(\phi)}, \sigma_{q(\phi)}^2, \lambda_{q(\phi)}; n, \sum_{i=1}^n e^{-x_i}, \mu_{\phi}, \sigma_{\phi}^2 \right) \right\} \quad (50)$$

and

$$\begin{aligned} f_{\text{Ex1}}^{\text{SN}}(x, y, z; a, b, c, d) &= \frac{1}{2} \log(y) - \sqrt{\frac{2}{\pi}} \int_{-\infty}^{\infty} \log\{\Phi(zt)\} \Phi(zt) e^{-t^2/2} dt \\ &\quad + a \left\{ x + z \sqrt{\frac{2y}{\pi(1+z^2)}} \right\} - 2b \exp\left(x + \frac{1}{2}y\right) \Phi\left(\frac{z\sqrt{y}}{\sqrt{z^2+1}}\right) \\ &\quad - \frac{1}{2d} \left[\left\{ x + z \sqrt{\frac{2y}{\pi(1+z^2)}} - c \right\}^2 + y \left\{ 1 - \frac{2z^2}{\pi(1+z^2)} \right\} \right]. \end{aligned}$$

Optimization problem (50) is considerably more challenging than its Normal counterpart. In particular, evaluations of the objective function and its derivatives require numerical integration.

We solved (50) for three Gumbel random samples of size $n = 5, 10$ and 20 and with $\sum_{i=1}^n e^{-x_i} \approx n$, corresponding to the mean of this sufficient statistic. The intractable integral in $f_{\text{Ex1}}^{\text{SN}}$ was approximated using a trapezoidal quadrature scheme similar to that described in Appendix B.2 of Wand et al. (2011). The limits of the trapezoidal grid were increased until the ratio of the global maximum and minimum absolute values of the integrand fell below 10^{-20} . The number of grid points was then doubled until the relative difference between two successive iterations was less than 10^{-20} . Multiple start locations and simulated annealing were used to locate global optima. Natural fixed-point iteration no

longer applies in this non-exponential family example and optimization of $f_{\text{Ex1}}^{\text{SN}}$ was accomplished using the *Broyden-Fletcher-Goldfarb-Shanno* quasi-Newton method via the `optim()` function in the R computing environment (R Development Core Team, 2016).

Figure 7 shows the optimal Skew-Normal q -density functions, together with the exact posterior density functions and those based on the Normal q -density restriction. We see that the Normal approximation is inferior for very low sample sizes, but that the approximations are about the same for moderate to large sample sizes.

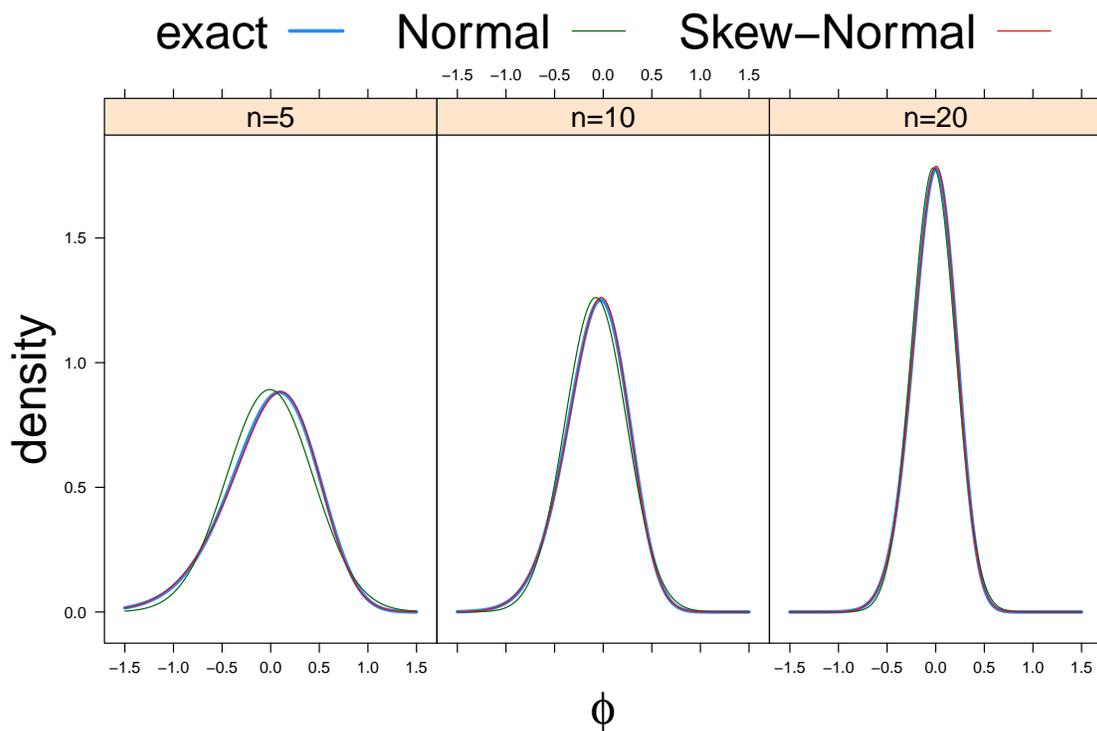


Figure 7: *Skew-Normal minimum Kullback-Leibler approximate posterior density functions for samples of size $n = 5, 10$ and 20 for the Example 1 Gumbel random sample setting. The exact posterior density functions and those based on restriction to the Normal family are also shown.*

6. Closing Remarks

We have taken a broad view of mean field variational Bayes with parametric pre-specification of one of the q -density components and coined the term ‘semiparametric mean field variational Bayes’ for this general approach. As well as laying out the general principles of

semiparametric mean field variational Bayes, we have provided an overview of the numerical issues attached to this methodology. Natural fixed-point iteration has been identified as a promising general approach to dealing with the Kullback-Leibler optimization problem and its attractive Riemannian gradient properties have been elucidated. Proof of convergence of a particular semiparametric mean field variational Bayes strategy appears to be too difficult a goal. However, for fixed-point iteration strategies, the spectral radius of the derivative matrix of the fixed-point update function is a reasonable diagnostic measure for checking convergence.

Acknowledgments

This research was partially supported by Australian Research Council Discovery Project DP110100061. We are grateful to Frances Kuo, Christian Wolff and Rob Womersley for their advice on aspects of this research. We also thank the editor and referees for their helpful comments.

Appendix A. Definitions and Derivations

A.1 Matrix definitions and identity

If \mathbf{A} is a $d \times d$ matrix then $\text{vec}(\mathbf{A})$ is the $d^2 \times 1$ vector obtained by stacking the columns of \mathbf{A} underneath each other in order from left to right. The inverse vec operator is denoted by vec^{-1} . In addition we let $\text{vech}(\mathbf{A})$ denote the $\frac{1}{2}d(d+1) \times 1$ vector obtained from $\text{vec}(\mathbf{A})$ by eliminating the above-diagonal entries of \mathbf{A} . If \mathbf{A} is symmetric then $\text{vech}(\mathbf{A})$ contains all of the unique entries of \mathbf{A} .

The derivations also require the commutation and duplication matrix notation of Magnus and Neudecker (1999). If \mathbf{A} is an arbitrary $d \times d$ matrix then the *commutation matrix of order d* , denoted by \mathbf{K}_d , is the $d^2 \times d^2$ matrix of zeroes and ones for which

$$\mathbf{K}_d \text{vec}(\mathbf{A}) = \text{vec}(\mathbf{A}^T).$$

If \mathbf{B} is a symmetric but otherwise arbitrary $d \times d$ matrix then the *duplication matrix of order d* is the $d^2 \times \frac{1}{2}d(d+1)$ matrix of zeroes and ones for which

$$\mathbf{D}_d \text{vech}(\mathbf{B}) = \text{vec}(\mathbf{B}).$$

The Moore-Penrose inverse of \mathbf{D}_d is

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

Note that

$$\mathbf{D}_d^+ \text{vec}(\mathbf{B}) = \text{vech}(\mathbf{B}). \quad (51)$$

Another useful notation is

$$\mathbf{Q}(\mathbf{A}) \equiv (\mathbf{A} \otimes \mathbf{1}^T) \odot (\mathbf{1}^T \otimes \mathbf{A})$$

for a general $m \times n$ matrix \mathbf{A} and $\mathbf{1}$ a $n \times 1$ vector of ones. The symbol \odot denotes element-wise product.

The following well-known matrix identity is used several times in the derivations:

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}(\mathbf{B}). \quad (52)$$

A.2 Derivative Matrix and Hessian Matrix Notation

Our summary of derivative-based optimization, and subsequent discussion, benefits from derivative vector and Hessian matrix notation. Such notation is not universal, and throughout this article we follow the conventions of Magnus and Neudecker (1999).

If \mathbf{h} is a \mathbb{R}^p -valued with argument $\mathbf{x} \in \mathbb{R}^d$ then the *derivative matrix* of \mathbf{h} with respect to \mathbf{x} , denoted by $\mathbf{D}_x \mathbf{h}(\mathbf{x})$, is the $p \times d$ matrix with (i, j) entry

$$\frac{\partial \mathbf{h}(\mathbf{x})_i}{\partial x_j}$$

A concrete derivative vector example is given in Section 2.3 of Wand (2014).

In the case $p = 1$, the *Hessian matrix* of \mathbf{h} with respect to \mathbf{x} is the $d \times d$ matrix

$$\mathbf{H}_x \mathbf{h}(\mathbf{x}) \equiv \mathbf{D}_x [\{\mathbf{D}_x \mathbf{h}(\mathbf{x})\}^T].$$

A.3 Example 2 Derivational Details

Here provide derivational details pertaining to Example 2 discussed in Section 2.2.3.

According to product restriction (22), the optimal q -density functions satisfy

$$\begin{aligned} q^*(\boldsymbol{\beta}, \mathbf{u}) &\propto \exp[E_{q(\sigma^2, a)} \log\{p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma^2, a)\}], \\ q^*(\sigma^2) &\propto \exp[E_{q(\boldsymbol{\beta}, \mathbf{u}, a)} \log\{p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma^2, a)\}] \\ \text{and } q^*(a) &\propto \exp[E_{q(\boldsymbol{\beta}, \mathbf{u}, \sigma^2)} \log\{p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma^2, a)\}] \end{aligned}$$

(e.g. Bishop, 2006, Section 10.1.1). Simple algebraic steps lead to the forms given in (23).

First we consider general pre-specified q -density families of the form $q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$, $\boldsymbol{\xi} \in \Xi$. With the help of (10) each of the terms in (27), can be expressed as follows:

$$\begin{aligned}
 \text{Entropy}\{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})\} &= - \int_{\mathbb{R}^{K+2}} \log\{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})\} q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi}) d\boldsymbol{\beta} d\mathbf{u}, \\
 \text{Entropy}\{q(\sigma^2)\} &= \log(B_{q(\sigma^2)}) + \frac{1}{2}(K+1) + \log\{\Gamma(\frac{1}{2}(K+1))\} \\
 &\quad - \frac{1}{2}(K+3)\text{digamma}\{\frac{1}{2}(K+1)\}, \\
 \text{Entropy}\{q(a)\} &= \log(B_{q(a)}) + 1 - 2\text{digamma}(1), \\
 E_q\{\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u})\} &= \mathbf{y}^T \{ \mathbf{X} E_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\boldsymbol{\beta}) + \mathbf{Z} E_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u}) \} \\
 &\quad - \mathbf{1}^T E_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})} \{ \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}) \} - \mathbf{1}^T \log(\mathbf{y}!), \\
 E_q\{\log p(\boldsymbol{\beta}, \mathbf{u} | \sigma^2)\} &= -\frac{1}{2}(p+K) \log(2\pi) - \frac{1}{2} p \log(\sigma_\beta^2) \\
 &\quad - \frac{1}{2} K \{ \log\{B_{q(\sigma^2)}\} - \text{digamma}\{\frac{1}{2}(K+1)\} \} \\
 &\quad - \frac{1}{2\sigma_\beta^2} \left[\| E_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\boldsymbol{\beta}) \|^2 + \text{tr}\{\text{Cov}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\boldsymbol{\beta})\} \right] \\
 &\quad - \frac{1}{2} \mu_{q(1/\sigma^2)} \left[\| E_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u}) \|^2 + \text{tr}\{\text{Cov}_{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})}(\mathbf{u})\} \right], \\
 E_q\{\log p(\sigma^2 | a)\} &= -\frac{1}{2} \log(\pi) - \frac{1}{2} \{ \log\{B_{q(a)}\} - \text{digamma}(1) \} \\
 &\quad - \frac{3}{2} \{ \log\{B_{q(\sigma^2)}\} - \text{digamma}\{\frac{1}{2}(K+1)\} \} \\
 &\quad - \mu_{q(1/a)} \mu_{q(1/\sigma^2)} \\
 \text{and } E_q\{\log p(a)\} &= -\frac{1}{2} \log(\pi) - \log(A) - \frac{3}{2} \{ \log\{B_{q(a)}\} - \text{digamma}(1) \} \\
 &\quad - \mu_{q(1/a)} / A^2.
 \end{aligned} \tag{53}$$

The $(\boldsymbol{\beta}, \mathbf{u})$ -localized approximate marginal log-likelihood expression given by (28) follows immediately from the relevant terms in (53).

If $q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\xi})$ is specified to be the $N(\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ density function then the terms in (27) that depend on $\boldsymbol{\xi} = (\boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$ are

$$\begin{aligned}
 \text{Entropy}\{q(\boldsymbol{\beta}, \mathbf{u}; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})\} &= \frac{1}{2}(p+K) \{1 + \log(2\pi)\} + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})}|, \\
 E_q\{\log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u})\} &= \mathbf{y}^T \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})} - \mathbf{1}^T \log(\mathbf{y}!) \\
 &\quad - \mathbf{1}^T \exp\{ \mathbf{C} \boldsymbol{\mu}_{q(\boldsymbol{\beta})} + \frac{1}{2} \text{diagonal}(\mathbf{C} \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})} \mathbf{C}^T) \} \\
 \text{and } E_q\{\log p(\boldsymbol{\beta}, \mathbf{u} | \sigma^2)\} &= -\frac{1}{2}(p+K) \log(2\pi) - \frac{1}{2} p \log(\sigma_\beta^2) \\
 &\quad - \frac{1}{2} K \{ \log\{B_{q(\sigma^2)}\} - \text{digamma}\{\frac{1}{2}(K+1)\} \} \\
 &\quad - \frac{1}{2\sigma_\beta^2} \left\{ \| \boldsymbol{\mu}_{q(\boldsymbol{\beta})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}) \right\} \\
 &\quad - \frac{1}{2} \mu_{q(1/\sigma^2)} \left\{ \| \boldsymbol{\mu}_{q(\boldsymbol{\mu})} \|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\mu})}) \right\}
 \end{aligned} \tag{54}$$

where $\mathbf{C} \equiv [\mathbf{X} \ \mathbf{Z}]$. The $(\boldsymbol{\beta}, \mathbf{u})$ -localized approximate marginal log-likelihood expression given by (29) follows immediately. An explicit expression for $\log \underline{p}(q; \boldsymbol{\mu}_{q(\boldsymbol{\beta}, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta}, \mathbf{u})})$, for use as a stopping criterion, can be formed by combining the relevant terms from (53) with those in (54).

A.4 Proof of (34)

In this proof, all appearances of \mathbf{D} and \mathbf{H} are assumed to be with respect to \mathbf{x} . Let

$$\mathbf{g}_{\text{NR}}(\mathbf{x}) \equiv \mathbf{x} - \{\mathbf{H} f(\mathbf{x})\}^{-1} \mathbf{D} f(\mathbf{x})^T.$$

Then, using (52),

$$\begin{aligned} d\mathbf{g}_{\text{NR}}(\mathbf{x}) &= d\mathbf{x} + \{\mathbf{H} f(\mathbf{x})\}^{-1} \{d\mathbf{H} f(\mathbf{x})\} \{\mathbf{H} f(\mathbf{x})\}^{-1} \mathbf{D} f(\mathbf{x})^T \\ &\quad - \{\mathbf{H} f(\mathbf{x})\}^{-1} d\mathbf{D} f(\mathbf{x})^T \\ &= d\mathbf{x} + \{\mathbf{H} f(\mathbf{x})\}^{-1} \text{vec}^{-1}[\mathbf{D} \text{vec}\{\mathbf{H} f(\mathbf{x})\} d\mathbf{x}] \{\mathbf{H} f(\mathbf{x})\}^{-1} \mathbf{D} f(\mathbf{x})^T \\ &\quad - \{\mathbf{H} f(\mathbf{x})\}^{-1} \mathbf{H} f(\mathbf{x}) d\mathbf{x} \\ &= \{\mathbf{H} f(\mathbf{x})\}^{-1} \text{vec} \left(\mathbf{I} \text{vec}^{-1}[\mathbf{D} \text{vec}\{\mathbf{H} f(\mathbf{x})\} d\mathbf{x}] \{\mathbf{H} f(\mathbf{x})\}^{-1} \mathbf{D} f(\mathbf{x})^T \right) \\ &= \{\mathbf{H} f(\mathbf{x})\}^{-1} \left([\mathbf{D} f(\mathbf{x}) \{\mathbf{H} f(\mathbf{x})\}^{-1}] \otimes \mathbf{I} \right) \mathbf{D} \text{vec}\{\mathbf{H} f(\mathbf{x})\} d\mathbf{x}. \end{aligned}$$

Therefore

$$\mathbf{D} \mathbf{g}_{\text{NR}}(\mathbf{x}) = \{\mathbf{H} f(\mathbf{x})\}^{-1} \left([\mathbf{D} f(\mathbf{x}) \{\mathbf{H} f(\mathbf{x})\}^{-1}] \otimes \mathbf{I} \right) \mathbf{D} \text{vec}\{\mathbf{H} f(\mathbf{x})\}.$$

Since $\mathbf{D} f(\mathbf{x}^*) = \mathbf{0}$, we get

$$\mathbf{D} \mathbf{g}_{\text{NR}}(\mathbf{x}^*) = \mathbf{O},$$

where \mathbf{O} is the $d \times d$ matrix with all entries equal to zero, and (34) follows immediately.

A.5 Lemmas and Proofs Required for Results 1 and 2

Lemma 1 *If*

$$q(\mathbf{x}; \boldsymbol{\eta}) = \exp\{\mathbf{T}(\mathbf{x})^T \boldsymbol{\eta} - A(\boldsymbol{\eta})\} h(\mathbf{x})$$

is an exponential family density function then

$$\mathbf{D}_{\boldsymbol{\eta}} \text{Entropy}\{q(\mathbf{x}; \boldsymbol{\eta})\} = -\boldsymbol{\eta}^T \mathbf{H}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}).$$

Proof of Lemma 1

Since

$$\text{Entropy}\{q(\mathbf{x}; \boldsymbol{\eta})\} = A(\boldsymbol{\eta}) - \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) \boldsymbol{\eta} - E[\log\{h(\mathbf{x})\}]$$

we then have

$$\mathbf{D}_{\boldsymbol{\eta}} \text{Entropy}\{q(\mathbf{x}; \boldsymbol{\eta})\} = \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) - \mathbf{D}_{\boldsymbol{\eta}} \{\mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) \boldsymbol{\eta}\} = \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) - \mathbf{D}_{\boldsymbol{\eta}} \{\boldsymbol{\eta}^T \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T\}.$$

Next,

$$\begin{aligned} d\{\boldsymbol{\eta}^T \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T\} &= (d\boldsymbol{\eta})^T \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T + \boldsymbol{\eta}^T d\{\mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T\} \\ &= \mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) d\boldsymbol{\eta} + \boldsymbol{\eta}^T \mathbf{D}_{\boldsymbol{\eta}} \{\mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})^T\} d\boldsymbol{\eta} \\ &= \{\mathbf{D}_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) + \boldsymbol{\eta}^T \mathbf{H}_{\boldsymbol{\eta}} A(\boldsymbol{\eta})\} d\boldsymbol{\eta} \end{aligned}$$

and so

$$D_{\boldsymbol{\eta}}\{\boldsymbol{\eta}^T D_{\boldsymbol{\eta}}A(\boldsymbol{\eta})^T\} = D_{\boldsymbol{\eta}}A(\boldsymbol{\eta}) + \boldsymbol{\eta}^T H_{\boldsymbol{\eta}}A(\boldsymbol{\eta}).$$

Hence,

$$D_{\boldsymbol{\eta}} \text{Entropy}\{q(\mathbf{x}; \boldsymbol{\eta})\} = D_{\boldsymbol{\eta}}A(\boldsymbol{\eta}) - \{D_{\boldsymbol{\eta}}A(\boldsymbol{\eta}) + \boldsymbol{\eta}^T H_{\boldsymbol{\eta}}A(\boldsymbol{\eta})\} = -\boldsymbol{\eta}^T H_{\boldsymbol{\eta}}A(\boldsymbol{\eta}).$$

Lemma 2 *Let s be a differentiable scalar-valued function of $\mathbf{x} \in \mathbb{R}^d$ and let $\mathbf{u} \in \mathbb{R}^k$ be one-to-one transformation of \mathbf{x} . Then*

$$D_{\mathbf{x}} s = (D_{\mathbf{u}} s) (D_{\mathbf{x}} \mathbf{u}).$$

Proof of Lemma 2

Lemma 2 is a restatement of Theorem 8, Chapter 5, of Magnus and Neudecker (1999).

Lemma 3 *Let $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ have a d -dimensional Multivariate Normal distribution. The natural statistic is*

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \text{vech}(\mathbf{x}\mathbf{x}^T) \end{bmatrix}$$

and corresponding mean parameter is $\boldsymbol{\tau} \equiv E\{\mathbf{T}(\mathbf{x})\}$. Then

$$D_{\boldsymbol{\tau}} \begin{bmatrix} \boldsymbol{\mu} \\ \text{vec}(\boldsymbol{\Sigma}) \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -(\mathbf{I} + \mathbf{K}_d)(\boldsymbol{\mu} \otimes \mathbf{I}) & \mathbf{D}_d \end{bmatrix}.$$

Proof of Lemma 3

The transformation from the common parameters to the mean parameters is

$$\boldsymbol{\tau} \equiv \begin{bmatrix} \boldsymbol{\tau}_1 \\ \boldsymbol{\tau}_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \text{vech}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T) \end{bmatrix}$$

and the inverse transformation is easily shown to be

$$\begin{bmatrix} \boldsymbol{\mu} \\ \text{vec}(\boldsymbol{\Sigma}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\tau}_1 \\ \mathbf{D}_d \boldsymbol{\tau}_2 - \text{vec}(\boldsymbol{\tau}_1 \boldsymbol{\tau}_1^T) \end{bmatrix}.$$

Hence

$$D_{\boldsymbol{\tau}} \begin{bmatrix} \boldsymbol{\mu} \\ \text{vec}(\boldsymbol{\Sigma}) \end{bmatrix} = \begin{bmatrix} D_{\boldsymbol{\tau}_1} \boldsymbol{\mu} & D_{\boldsymbol{\tau}_2} \boldsymbol{\mu} \\ D_{\boldsymbol{\tau}_1} \text{vec}(\boldsymbol{\Sigma}) & D_{\boldsymbol{\tau}_2} \text{vec}(\boldsymbol{\Sigma}) \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -D_{\boldsymbol{\tau}_1} \text{vec}(\boldsymbol{\tau}_1 \boldsymbol{\tau}_1^T) & \mathbf{D}_d \end{bmatrix}.$$

To obtain an explicit expression for the bottom left-hand block we note that

$$d\text{vec}(\boldsymbol{\tau}_1 \boldsymbol{\tau}_1^T) = (\mathbf{I} + \mathbf{K}_d) \text{vec}\{(d\boldsymbol{\tau}_1) \boldsymbol{\tau}_1^T\}.$$

Then, with the help of (52),

$$\text{vec}\{(d\boldsymbol{\tau}_1) \boldsymbol{\tau}_1^T\} = \text{vec}\{\mathbf{I}(d\boldsymbol{\tau}_1) \boldsymbol{\tau}_1^T\} = (\boldsymbol{\tau}_1 \otimes \mathbf{I}) d\boldsymbol{\tau}_1.$$

Hence,

$$D_{\tau_1} \text{vec}(\Sigma) = (\mathbf{I} + \mathbf{K}_d)(\tau_1 \otimes \mathbf{I}) = (\mathbf{I} + \mathbf{K}_d)(\boldsymbol{\mu} \otimes \mathbf{I})$$

and the lemma follows immediately.

Lemma 4 *Let*

$$\phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \equiv (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$$

denote the d -variate $N(\boldsymbol{\mu}, \Sigma)$ density function. Then

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

Proof of Lemma 4

First note that

$$(2\pi)^{d/2} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = |\Sigma|^{-1/2} \exp[-\frac{1}{2} \text{tr}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}\}].$$

Then, using the identity $\text{tr}(\mathbf{A}^T \mathbf{B}) = \text{vec}(\mathbf{A})^T \text{vec}(\mathbf{B})$,

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

Therefore, by Theorem 6, Chapter 5, of Magnus and Neudecker (1999),

$$D_{\text{vec}(\Sigma)} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{2} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \text{vec}[\Sigma^{-1}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} - \mathbf{I}\}]^T.$$

Also,

$$\begin{aligned} |\Sigma|^{1/2} (2\pi)^{d/2} d_{\boldsymbol{\mu}} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) &= \exp[-\frac{1}{2} \text{tr}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}\}] \\ &\quad \times [-\frac{1}{2} \text{tr}\{d_{\boldsymbol{\mu}}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\} \Sigma^{-1}\}] \\ &= |\Sigma|^{1/2} (2\pi)^{d/2} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \text{tr}\{(\mathbf{x} - \boldsymbol{\mu})(d_{\boldsymbol{\mu}})^T \Sigma^{-1}\} \\ &= |\Sigma|^{1/2} (2\pi)^{d/2} \phi(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})^T d\boldsymbol{\mu} \end{aligned}$$

which simplifies to

$$d_{\boldsymbol{\mu}}\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})^T d\boldsymbol{\mu}.$$

The second differential with respect to $\boldsymbol{\mu}$ is then

$$\begin{aligned} d_{\boldsymbol{\mu}}^2\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \{d_{\boldsymbol{\mu}}\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})^T d\boldsymbol{\mu} + \phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(-d\boldsymbol{\mu})^T d\boldsymbol{\mu} \\ &= \{\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})^T d\boldsymbol{\mu}\}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})^T d\boldsymbol{\mu} + \phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1}(-d\boldsymbol{\mu})^T d\boldsymbol{\mu} \\ &= (d\boldsymbol{\mu})^T \left(\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})[\boldsymbol{\Sigma}^{-1}\{(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{I}\}] \right) d\boldsymbol{\mu}. \end{aligned}$$

Hence, using Theorem 6, Chapter 6, of Magnus and Neudecker (1999)

$$H_{\boldsymbol{\mu}}\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})[\boldsymbol{\Sigma}^{-1}\{(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{I}\}] = 2\text{vec}^{-1}\left\{D_{\text{vec}(\boldsymbol{\Sigma})}\phi(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})^T\right\}.$$

A.6 Derivation of Result 2

To make the derivation less cumbersome we will suppress the subscripts on the mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. As in Wand (2014) we work with the natural statistic and natural parameter pair

$$\boldsymbol{T}(\boldsymbol{\phi}) = \begin{bmatrix} \boldsymbol{\phi} \\ \text{vech}(\boldsymbol{\phi}\boldsymbol{\phi}^T) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\boldsymbol{D}_d\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}.$$

The mean parameter vector is

$$\boldsymbol{\tau} = E\{\boldsymbol{T}(\boldsymbol{\phi})\} = \begin{bmatrix} \boldsymbol{\mu} \\ \text{vech}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T) \end{bmatrix}.$$

In Lemma 3 in Appendix A.5 we show that

$$D_{\boldsymbol{\tau}} \begin{bmatrix} \boldsymbol{\mu} \\ \text{vec}(\boldsymbol{\Sigma}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{I} & \mathbf{0} \\ -(\boldsymbol{I} + \boldsymbol{K}_d)(\boldsymbol{\mu} \otimes \boldsymbol{I}) & \boldsymbol{D}_d \end{bmatrix}$$

and so Result 1(d) becomes

$$\begin{aligned} &\begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\boldsymbol{D}_d^T\text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} \\ &= \begin{bmatrix} \boldsymbol{I} & -(\boldsymbol{\mu}^T \otimes \boldsymbol{I})(\boldsymbol{I} + \boldsymbol{K}_d) \\ \mathbf{0} & \boldsymbol{D}_d^T \end{bmatrix} \begin{bmatrix} D_{\boldsymbol{\mu}}\text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\ [D_{\text{vec}(\boldsymbol{\Sigma})}\text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \end{bmatrix} \end{aligned} \tag{55}$$

where we have used the fact that $\boldsymbol{K}_d^T = \boldsymbol{K}_d$. Using (52), and the fact that

$$\text{vec}^{-1}[D_{\text{vec}(\boldsymbol{\Sigma})}\text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T]$$

is symmetric, the first component of (55) is equivalent to

$$\begin{aligned}
 \Sigma^{-1} \boldsymbol{\mu} &= \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &\quad - (\boldsymbol{\mu}^T \otimes \mathbf{I})(\mathbf{I} + \mathbf{K}_d) \mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &= \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &\quad - (\boldsymbol{\mu}^T \otimes \mathbf{I})(\mathbf{I} + \mathbf{K}_d) \text{vec} \left(\text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \right) \\
 &= \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &\quad - 2(\boldsymbol{\mu}^T \otimes \mathbf{I}) \text{vec} \left(\text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \right) \\
 &= \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &\quad - 2 \text{vec} \left(\text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \boldsymbol{\mu} \right) \\
 &= \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 &\quad - 2 \text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \boldsymbol{\mu}.
 \end{aligned} \tag{56}$$

The second component (55) is equivalent to

$$-\frac{1}{2} \mathbf{D}_d^T \text{vec}(\boldsymbol{\Sigma}) = \mathbf{D}_d^T [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T]^T$$

which, under the constraint that $\boldsymbol{\Sigma}$ is symmetric, is equivalent to

$$\boldsymbol{\Sigma}^{-1} = -2 \text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T]. \tag{57}$$

In view of relationships (56) and (57), the natural fixed-point iteration scheme becomes

$$\left\{ \begin{array}{l}
 \boldsymbol{\Sigma}_{\text{new}}^{-1} \boldsymbol{\mu}_{\text{new}} = [\mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}]_{\boldsymbol{\mu}=\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{\text{old}}}^T \\
 \quad - 2 \text{vec}^{-1} \left([\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}]_{\boldsymbol{\mu}=\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{\text{old}}}^T \right) \boldsymbol{\mu}_{\text{old}} \\
 \boldsymbol{\Sigma}_{\text{new}} = \left\{ -2 \text{vec}^{-1} \left([\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}]_{\boldsymbol{\mu}=\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{\text{old}}}^T \right) \right\}^{-1}
 \end{array} \right.$$

where $(\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}_{\text{old}})$ and $(\boldsymbol{\mu}_{\text{new}}, \boldsymbol{\Sigma}_{\text{new}})$, respectively, denote the old and new values of $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. The following simplification ensues:

$$\left\{ \begin{array}{l}
 \boldsymbol{\mu}_{\text{new}} = \boldsymbol{\mu}_{\text{old}} + \boldsymbol{\Sigma}_{\text{new}} [\mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}]_{\boldsymbol{\mu}=\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{\text{old}}}^T \\
 \boldsymbol{\Sigma}_{\text{new}} = \left\{ -2 \text{vec}^{-1} \left([\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}]_{\boldsymbol{\mu}=\boldsymbol{\mu}_{\text{old}}, \boldsymbol{\Sigma}=\boldsymbol{\Sigma}_{\text{old}}}^T \right) \right\}^{-1}
 \end{array} \right.$$

which is equivalent to the following updating scheme:

$$\left\{ \begin{array}{l}
 \mathbf{v} \leftarrow \mathbf{D}_{\boldsymbol{\mu}} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T \\
 \boldsymbol{\Sigma} \leftarrow \left(-2 \text{vec}^{-1} [\mathbf{D}_{\text{vec}(\boldsymbol{\Sigma})} \text{NonEntropy}\{q(\boldsymbol{\phi}; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] \right)^{-1} \\
 \boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{v}
 \end{array} \right. \tag{58}$$

as given in Wand (2014). However, from Lemma 4 in Appendix A.5 (see also Appendix A of Opper and Archambeau, 2009) and the fact that $\text{NonEntropy}\{q(\phi; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}$ is an expectation with respect to the $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ density function we have

$$\text{vec}^{-1}[\text{D}_{\text{vec}(\boldsymbol{\Sigma})}\text{NonEntropy}\{q(\phi; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}^T] = \frac{1}{2}\mathbf{H}_{\boldsymbol{\mu}}\text{NonEntropy}\{q(\phi; \boldsymbol{\mu}, \boldsymbol{\Sigma})\}$$

which leads to a somewhat more elegant alternative to (58) given in Result 2.

A.7 Derivation of the Derivative Matrix of $\mathbf{g}_{\text{Ex}2}$

From Section 4.3, the fixed-point iteration updating function is of the form

$$\mathbf{g}_{\text{Ex}2} \left(\begin{bmatrix} \boldsymbol{\mu} \\ \text{vech}(\boldsymbol{\Sigma}) \end{bmatrix} \right) \equiv \begin{bmatrix} \mathbf{g}_{\text{Ex}2, \boldsymbol{\mu}} \\ \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) \end{bmatrix}$$

where

$$\begin{aligned} \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} &\equiv [\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}]^{-1}, \\ \mathbf{g}_{\text{Ex}2, \boldsymbol{\mu}} &\equiv \boldsymbol{\mu} + \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} [\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}] \end{aligned}$$

$$\text{and } \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \exp\{\mathbf{C}\boldsymbol{\mu} + \frac{1}{2}\text{diagonal}(\mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T)\}.$$

Note that the dependence of $\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}$ and $\mathbf{g}_{\text{Ex}2}$ on \mathbf{y} , \mathbf{C} and \mathbf{M} is suppressed here.

The derivative matrix of $\mathbf{g}_{\text{Ex}2}$ with respect to $[\boldsymbol{\mu} \text{vech}(\boldsymbol{\Sigma})]^T$ is

$$\text{D} \begin{bmatrix} \boldsymbol{\mu} \\ \text{vech}(\boldsymbol{\Sigma}) \end{bmatrix} \mathbf{g}_{\text{Ex}2} = \begin{bmatrix} \text{D}_{\boldsymbol{\mu}} \mathbf{g}_{\text{Ex}2, \boldsymbol{\mu}} & \text{D}_{\text{vech}(\boldsymbol{\Sigma})} \mathbf{g}_{\text{Ex}2, \boldsymbol{\mu}} \\ \text{D}_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) & \text{D}_{\text{vech}(\boldsymbol{\Sigma})} \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) \end{bmatrix}.$$

We now give explicit expressions for each of these four components of the derivative matrix. It is more efficient, notationally, to first obtain expressions for the derivatives of $\text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}})$.

A.7.1 EXPRESSION FOR $\text{D}_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}})$

$$\text{D}_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) = -\mathbf{D}_{p+K}^+ (\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} \otimes \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) \mathbf{Q}(\mathbf{C})^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} \quad (59)$$

Derivation:

Using the second rule in Section 3.3.5 of Wand (2002), (52) and (51),

$$\begin{aligned} d_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}}) &= \mathbf{D}_{p+K}^+ d_{\boldsymbol{\mu}} \text{vec} \left([\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}]^{-1} \right) \\ &= -\mathbf{D}_{p+K}^+ \text{vec} \left\{ \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} \left(d_{\boldsymbol{\mu}} [\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}] \right) \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} \right\} \\ &= -\mathbf{D}_{p+K}^+ \left(\mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} \otimes \mathbf{G}_{\text{Ex}2, \boldsymbol{\Sigma}} \right) \text{vec} [\mathbf{C}^T \text{diag}\{d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C}]. \end{aligned}$$

From Theorem 2(b) of Wand (2014),

$$\text{vec}[\mathbf{C}^T \text{diag}\{d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C}] = \mathbf{Q}(\mathbf{C})^T d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

Lastly, we use the chain rule in Section 3.3.2 of Wand (2002) to get

$$d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = d_{\boldsymbol{\mu}} \exp\{\mathbf{C}\boldsymbol{\mu} + \frac{1}{2} \text{diagonal}(\mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T)\} = \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} d_{\boldsymbol{\mu}}.$$

Combining, we then obtain

$$d_{\boldsymbol{\mu}} \text{vec}(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}) = -\mathbf{D}_{p+K}^+ \left(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} \otimes \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} \right) \mathbf{Q}(\mathbf{C})^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} d_{\boldsymbol{\mu}}$$

and the stated expression then follow from Theorem 6, Chapter 5, of Magnus and Neudecker (1999).

A.7.2 EXPRESSION FOR $\mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})} \text{vech}(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}})$

$$\begin{aligned} \mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})} \text{vech}(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}) &= -\frac{1}{2} \mathbf{D}_{p+K}^+ (\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} \otimes \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}) \\ &\quad \times \mathbf{Q}(\mathbf{C})^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{Q}(\mathbf{C}) \mathbf{D}_{p+K}. \end{aligned} \quad (60)$$

Derivation:

The derivation is similar to that for $\mathbf{D}_{\boldsymbol{\mu}} \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}$. It differs in that it requires $d_{\boldsymbol{\Sigma}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ rather than $d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. This entails

$$d_{\boldsymbol{\Sigma}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = d_{\boldsymbol{\Sigma}} \exp\{\mathbf{C}\boldsymbol{\mu} + \frac{1}{2} \text{diagonal}(\mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T)\} = \frac{1}{2} \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} d_{\boldsymbol{\Sigma}} \text{diagonal}(\mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T).$$

But Theorem 2(a) of Wand (2014) gives

$$d_{\boldsymbol{\Sigma}} \text{diagonal}(\mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T) = \mathbf{Q}(\mathbf{C}) d_{\boldsymbol{\Sigma}} \text{vec}(\boldsymbol{\Sigma}) = \mathbf{Q}(\mathbf{C}) \mathbf{D}_{p+K} d_{\boldsymbol{\Sigma}} \text{vech}(\boldsymbol{\Sigma})$$

which leads to the stated result.

A.7.3 EXPRESSION FOR $\mathbf{D}_{\boldsymbol{\mu}} \mathbf{g}_{\text{Ex2}, \boldsymbol{\mu}}$

$$\mathbf{D}_{\boldsymbol{\mu}} \mathbf{g}_{\text{Ex2}, \boldsymbol{\mu}} = ([\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I}) \mathbf{D}_{p+K} \mathbf{D}_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}})$$

where $\mathbf{D}_{\boldsymbol{\mu}} \text{vech}(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}})$ is given by (59).

Derivation:

Using the second rule in Section 3.3.4 of Wand (2002) for differentiation of matrix products,

$$\begin{aligned} d_{\boldsymbol{\mu}} \mathbf{g}_{\text{Ex2}, \boldsymbol{\mu}} &= d_{\boldsymbol{\mu}} + d_{\boldsymbol{\mu}} \left(\mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} [\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}] \right) \\ &= d_{\boldsymbol{\mu}} + (d_{\boldsymbol{\mu}} \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}) [\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}] \\ &\quad - \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} [\mathbf{C}^T d_{\boldsymbol{\mu}} \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \mathbf{M} d_{\boldsymbol{\mu}}] \\ &= d_{\boldsymbol{\mu}} + \text{vec}\{ \mathbf{I} (d_{\boldsymbol{\mu}} \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}}) [\mathbf{C}^T \{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}] \} \\ &\quad - \mathbf{G}_{\text{Ex2}, \boldsymbol{\Sigma}} [\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} \mathbf{C} + \mathbf{M}] d_{\boldsymbol{\mu}}, \end{aligned}$$

where we have used the fact that $(d_{\boldsymbol{\mu}}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})[\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]$ is a column vector. Application of (52) leads to

$$\begin{aligned}
 d_{\boldsymbol{\mu}}\mathbf{g}_{\text{Ex2},\boldsymbol{\mu}} &= d_{\boldsymbol{\mu}} + ([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I})d_{\boldsymbol{\mu}}\text{vec}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}) \\
 &\quad - \mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}[\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\}\mathbf{C} + \mathbf{M}]d_{\boldsymbol{\mu}} \\
 &= \left\{ \mathbf{I} + ([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I})\mathbf{D}_{\boldsymbol{\mu}}\text{vec}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}) \right. \\
 &\quad \left. - \mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}^{-1} \right\} d_{\boldsymbol{\mu}} \\
 &= ([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I})\mathbf{D}_{p+K}\mathbf{D}_{\boldsymbol{\mu}}\text{vech}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})d_{\boldsymbol{\mu}}.
 \end{aligned}$$

The given expression follows from Theorem 6, Chapter 5, of Magnus and Neudecker (1999).

A.7.4 EXPRESSION FOR $\mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})}\mathbf{g}_{\text{Ex2},\boldsymbol{\mu}}$

$$\begin{aligned}
 \mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})}\mathbf{g}_{\text{Ex2},\boldsymbol{\mu}} &= \left([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I} \right) \mathbf{D}_{p+K}\mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})}\text{vech}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}) \\
 &\quad - \frac{1}{2}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\}\mathbf{Q}(\mathbf{C})\mathbf{D}_{p+K}
 \end{aligned}$$

where $\mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})}\text{vech}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})$ is given by (60).

Derivation:

Dealing with matrix products via the second rule in Section 3.3.4 of Wand (2002) we obtain

$$\begin{aligned}
 d_{\boldsymbol{\Sigma}}\mathbf{g}_{\text{Ex2},\boldsymbol{\mu}} &= d_{\boldsymbol{\Sigma}}\left(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}[\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]\right) \\
 &= (d_{\boldsymbol{\Sigma}}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})[\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}] \\
 &\quad - \mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{C}^T d_{\boldsymbol{\Sigma}}\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
 &= \text{vec}\left(\mathbf{I}(d_{\boldsymbol{\Sigma}}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})[\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]\right) \\
 &\quad - \frac{1}{2}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\}\mathbf{Q}(\mathbf{C})d\text{vec}(\boldsymbol{\Sigma}) \\
 &= ([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I})\mathbf{D}_{p+K}d\text{vech}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}) \\
 &\quad - \frac{1}{2}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\}\mathbf{Q}(\mathbf{C})\mathbf{D}_{p+K}d\text{vech}(\boldsymbol{\Sigma}) \\
 &= ([\mathbf{C}^T\{\mathbf{y} - \boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\} - \mathbf{M}\boldsymbol{\mu}]^T \otimes \mathbf{I}) \\
 &\quad \times \mathbf{D}_{p+K}\mathbf{D}_{\text{vech}(\boldsymbol{\Sigma})}\text{vech}(\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}})d\text{vech}(\boldsymbol{\Sigma}) \\
 &\quad - \frac{1}{2}\mathbf{G}_{\text{Ex2},\boldsymbol{\Sigma}}\mathbf{C}^T \text{diag}\{\boldsymbol{\omega}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\}\mathbf{Q}(\mathbf{C})\mathbf{D}_{p+K}d\text{vech}(\boldsymbol{\Sigma}).
 \end{aligned}$$

Once again we call upon Theorem 6, Chapter 5, of Magnus and Neudecker (1999) to complete the derivation.

References

- A.S. Ackleh, E.J. Allen, R.B. Hearfott, and P. Seshaiyer. *Classical and Modern Numerical Analysis: Theory, Methods and Practice*. Chapman & Hall, CRC Press, Boca Raton, Florida, 2010.
- S. Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10:251–276, 1998.
- C. Archambeau, D. Cornford, M. Opper, and J. Shawe-Taylor. Gaussian process approximations of stochastic differential equations. *Journal of Machine Learning Research: Workshop and Conference Proceedings*, 1:1–16, 2007.
- A. Azzalini and A. Dalla Valle. The Multivariate Skew-Normal distribution. *Biometrika*, 83:715–726, 1996.
- D. Barber and C.M. Bishop. Ensemble learning for multi-layer networks. In M.I. Jordan, K.J. Kearns, and S.A. Solla, editors, *Advances in Neural Information Processing Systems 10*, pages 395–401, 1997.
- M.J. Beal and Z. Ghahramani. Variational Bayesian learning of directed graphical models with hidden variables. *Bayesian Analysis*, 4:793–832, 2006.
- C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, 2006.
- L. Bottou. Stochastic learning. In O. Bousquet, U. von Luxburg, and G. Rätsch, editors, *Advanced Lectures on Machine Learning*, pages 146–168, 2004.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- E. Challis and D. Barber. Gaussian Kullback-Leibler approximate inference. *Journal of Machine Learning Research*, 14:2239–2286, 2013.
- Y.H. Dai and Y. Yuan. A nonlinear conjugate gradient method with a strong global convergence property. *SIAM Journal on Optimization*, 10:177–182, 1999.
- S.J. Gershman, M.D. Hoffman, and D.M. Blei. Nonparametric variational inference. In *Proceedings of the Twenty Ninth International Conference on Machine Learning*, pages 633–670, 2012.
- G.H. Givens and J.A. Hoeting. *Computational Statistics*. Wiley, Hoboken, New Jersey, 2005.
- M.R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of Research of the National Bureau of Standards*, 49:409–436, 1952.
- G.E. Hinton and D. van Camp. Keeping neural networks simple by minimizing description length of the weights. In L. Pitt, editor, *Proceedings of the Sixth Annual Conference on Computational Learning Theory*, pages 5–13, 1993.

- M.D. Hoffman, D.M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14:1303–1347, 2013.
- A. Honkela and H. Valpola. Unsupervised variational Bayesian learning of nonlinear models. *Advances in Neural Information Processing Systems*, 17:593–600, 2005.
- A. Honkela, H. Valpola, A. Ilin, and J. Karhunen. Blind separation of nonlinear mixtures by variational Bayesian learning. *Digital Signal Processing*, 17:914–934, 2007.
- A. Honkela, M. Tornio, T. Raiko, and J. Karhunen. Natural conjugate gradient in variational inference. In M. Ishikawa, K. Doya, H. Miyamoto, and T. Yamakawa, editors, *Proceedings of the Fourteenth International Conference on Neural Information Processing*, pages 305–314, 2008.
- A. Honkela, T. Raiko, M. Kuusela, M. Tornio, and J. Karhunen. Approximate Riemannian conjugate gradient learning for fixed-form variational Bayes. *Journal of Machine Learning Research*, 11:3235–3268, 2010.
- D.A. Knowles and T.P. Minka. Non-conjugate message passing for multinomial and binary regression. In J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, 2011.
- H. Lappalainen and A. Honkela. Bayesian nonlinear independent component analysis by multilayer perceptrons. In M. Girolami, editor, *Advances in Independent Component Analysis*, pages 93–121, 2000.
- D.G. Luenberger and Y. Ye. *Linear and Nonlinear Programming, Third Edition*. Springer, New York, 2008.
- J. Luts and M.P. Wand. Variational inference for count response semiparametric regression. *Bayesian Analysis*, 10:991–1023, 2015.
- J.R. Magnus and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics, Revised Edition*. Wiley, Chichester, UK, 1999.
- M. Menictas and M.P. Wand. Variational inference for heteroscedastic semiparametric regression. *Australian and New Zealand Journal of Statistics*, 57:119–138, 2015.
- M.K. Murray and J.W. Rice. *Differential Geometry and Statistics*. Chapman & Hall, London, 1993.
- J.A. Nelder and R. Mead. A simplex method for function minimization. *Computer Journal*, 7:308–313, 1965.
- S.E. Neville, J.T. Ormerod, and M.P. Wand. Mean field variational Bayes for continuous sparse signal shrinkage: pitfalls and remedies. *Electronic Journal of Statistics*, 8:1113–1151, 2014.
- H. Nickisch and C.E. Rasmussen. Approximations for binary Gaussian process classification. *Journal of Machine Learning Research*, 9:2035–2078, 2008.

- M. Opper and C. Archambeau. The variational Gaussian approximation revisited. *Neural Computation*, 21:786–792, 2009.
- J.T. Ormerod and M.P. Wand. Explaining variational approximations. *The American Statistician*, 64:140–153, 2010.
- J.M. Ortega. *Numerical Analysis: A Second Course*. SIAM, Philadelphia, 1990.
- T. Pham, J.T. Ormerod, and M.P. Wand. Mean field variational Bayesian inference for nonparametric regression with measurement error. *Computational Statistics and Data Analysis*, 68:375–387, 2013.
- W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery. *Numerical Recipes: The Art of Scientific Computing, Third Edition*. Cambridge University Press, New York, 2007.
- T. Raiko, H. Valpola, M. Harva, and J. Karhunen. Building blocks for variational Bayesian learning of latent variable models. *Journal of Machine Learning Research*, 8:155–201, 2007.
- R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016. URL <http://www.R-project.org/>.
- T. Salimans and D.A. Knowles. Fixed-form variational posterior approximation through stochastic linear regression. *Bayesian Analysis*, 8:837–882, 2013.
- M. Sato. Online model selection based on the variational Bayes. *Neural Computation*, 13:1649–1681, 2001.
- L.S.L. Tan and D.J. Nott. Variational inference for generalized linear mixed models using partially noncentered parametrizations. *Statistical Science*, 28:168–188, 2013.
- M.J. Wainwright and M.I. Jordan. Graphical models, exponential families and variational inference. *Foundations and Trends in Machine Learning*, 1:1–305, 2008.
- M.P. Wand. Vector differential calculus in statistics. *The American Statistician*, 56:55–62, 2002.
- M.P. Wand. Fully simplified Multivariate Normal updates in non-conjugate variational message passing. *Journal of Machine Learning Research*, 15:1351–1369, 2014.
- M.P. Wand, J.T. Ormerod, S.A. Padoan, and R. Frühwirth. Mean field variational Bayes for elaborate distributions. *Bayesian Analysis*, 6:847–900, 2011.
- B. Wang and D.M. Titterton. Inadequacy of interval estimates corresponding to variational Bayes approximations. In Z. Ghahramani and R. Cowell, editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, pages 373–380, 2005.
- O. Zobay. Variational Bayesian inference with Gaussian-mixture approximations. *Electronic Journal of Statistics*, 8:355–389, 2014.