



VARIATIONAL INFERENCE FOR HETEROSCEDASTIC SEMIPARAMETRIC REGRESSION

MARIANNE MENICTAS* AND MATT P. WAND

University of Technology Sydney

Summary

We develop fast mean field variational methodology for Bayesian heteroscedastic semi-parametric regression, in which both the mean and variance are smooth, but otherwise arbitrary, functions of the predictors. Our resulting algorithms are purely algebraic, devoid of numerical integration and Monte Carlo sampling. The locality property of mean field variational Bayes implies that the methodology also applies to larger models possessing variance function components. Simulation studies indicate good to excellent accuracy, and considerable time savings compared with Markov chain Monte Carlo. We also provide some illustrations from applications.

Key words: Approximate Bayesian inference; mean field variational Bayes; non-conjugate variational message passing; variance function estimation.

1. Introduction

Data sets that are big in terms of volume and/or velocity are becoming widespread and there is a strong imperative for the development of fast, flexible and extendable methodology for processing such data. Semiparametric regression (e.g. Ruppert, Wand & Carroll 2003, 2009) is an important class of flexible statistical models and methods, but is mainly geared towards moderate sample sizes and batch processing. Recently, Luts, Broderick & Wand (2014) developed new nonparametric regression methodology specifically tailored to high volume/velocity situations. The present article extends their general approach to accommodate possible heteroscedasticity. Whilst we focus on univariate and bivariate nonparametric regression and additive models, the modularity of our approach allows easy extension to more complex models. The accommodation of heteroscedasticity is aided by the variational inference technique known as non-conjugate variational message passing (e.g. Knowles & Minka 2011). The particular form of non-conjugate variational message passing that we use involves approximating particular posterior density functions by Multivariate Normal density functions, which is an established practice within the variational approximation literature (e.g. Hinton & van Camp 1993; Barber & Bishop 1997; Raiko *et al.* 2007; Challis & Barber 2013).

In nonparametric regression it is common to ignore heteroscedasticity in the data and invoke the constant variance assumption. Figure 1 contains nonparametric regression examples, based on the R function `smooth.spline()` (R Core Team 2015) with default settings.

*Author to whom correspondence should be addressed.

School of Mathematical Sciences, University of Technology Sydney, Broadway, 2007, Australia.
e-mail: marianne.menictas@uts.edu.au.

Acknowledgements. This research was partially supported by an Australian Postgraduate Award and Australian Research Council Discovery Project DP110100061. The authors are grateful to Jan Luts, an associate editor and two referees for their comments on this research.

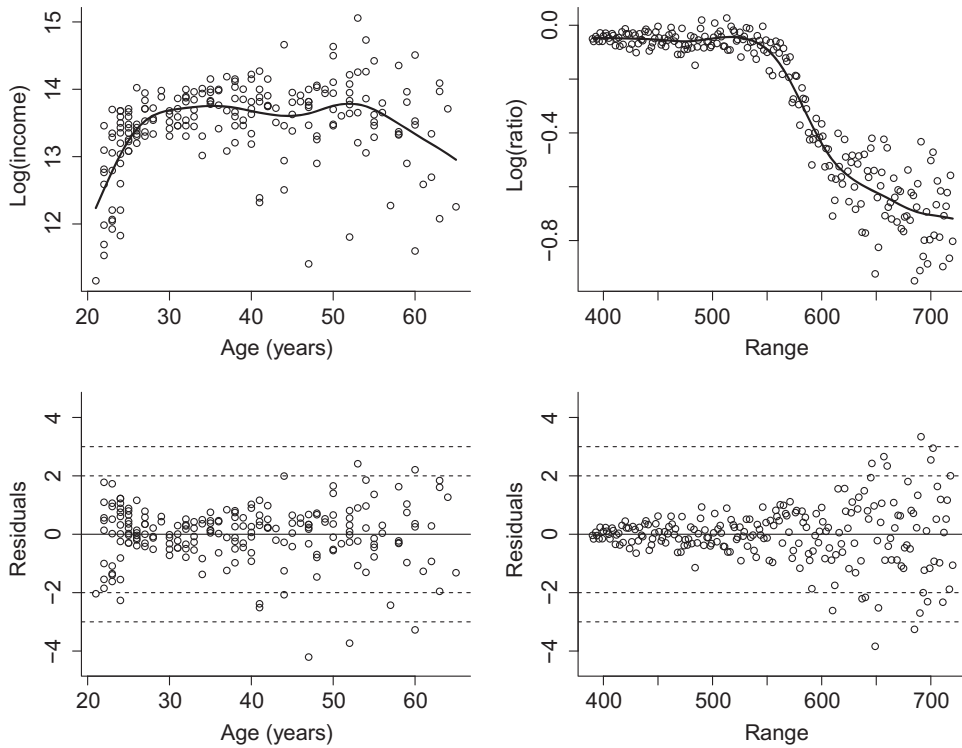


Figure 1. Two example nonparametric regression fits and corresponding standardized residual plots. Horizontal dashed lines at ± 2 and ± 3 aid assessment of standard normality of the standardized residuals.

Both data sets appear in Ruppert *et al.* (2003) and are first described in its Sections 5.3 and 2.7 respectively. The scatterplots and the standardized residual plots show that there is significant heteroscedasticity, which if ignored will lead to erroneous inference, for example incorrect prediction intervals.

Heteroscedastic nonparametric regression overcomes the problems apparent in Figure 1. For a set of regression data, (x_i, y_i) , $1 \leq i \leq n$, it involves replacement of the homoscedastic nonparametric model

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

by

$$E(y_i) = f(x_i), \quad \text{Var}(y_i) = g(x_i), \quad (1)$$

with the variance function g to be estimated simultaneously with the mean function f . Several approaches to fitting (1) now exist (e.g. Ruppert *et al.* 2003; Rigby & Stasinopoulos 2005; Crainiceanu *et al.* 2007). A particularly attractive approach, from an extendability standpoint, involves graphical model representations of mixed model-based penalized splines (Wand, 2009). This allows one to take advantage of the growing body of methodology and software for approximate inference in general graphical models. Graphical model-based Bayesian inference engines such as BUGS (Spiegelhalter *et al.*, 2003), Infer.NET (Minka

et al. 2014) and Stan (Stan Development Team 2013) now accommodate a wide range of nonparametric and semiparametric regression models (e.g. Marley & Wand 2010; Luts *et al.* 2015). However, the form of the approximate inference differs markedly among the various inference engines. The engines BUGS and Stan each use relatively slow but highly accurate Markov Chain Monte Carlo (MCMC) methodology whereas Infer.NET uses faster, but less accurate, deterministic approximations such as mean field variational Bayes (MFVB). The latter type of methodology is our focus here. A relatively new modification of MFVB, non-conjugate variational message passing, is shown to be particularly useful for handling heteroscedasticity since it results in algorithms that involve only closed form algebraic expressions.

The modularity of the graphical model and MFVB approaches implies that methodology for handling heteroscedasticity in semiparametric regression applies to arbitrarily large models. We call this the *locality* property of MFVB and it is described in section 3 of Wand *et al.* (2011). If a very large graphical model includes a component where one variable is modelled as a heteroscedastic nonparametric regression function of another variable then variational inference for the parameters in that part of the graph can be achieved using the methodology developed here.

Variational methodology for heteroscedastic regression models has also been developed by Lázaro-Gredilla & Titsias (2011) and Nott, Tran & Leng (2012), although the latter article was confined to linear mean and log-variance functions. Lazaro-Gredilla & Titsias (2011) achieved simultaneous nonparametric mean and variance function estimation via an elegant Gaussian process approach. The ensuing nonparametric function estimators are full-rank and their variance function estimation strategy relies on Gauss-Hermite quadrature. Our approach differs by using low-rank penalized splines and a non-conjugate MFVB strategy that provides closed form updates, making it more amenable to high volume/velocity data.

Section 2 gives a description of the Bayesian penalized spline model for simultaneous mean and variance function estimation based on univariate data. The variational inference methodology is provided in Section 3 and our main algorithm is presented there. Section 4 provides numerical illustrations which give evidence of the speed achieved by non-conjugate MFVB. In Sections 5–7 we describe extensions to bivariate nonparametric regression, additive models and real-time semiparametric regression. Some concluding remarks are made in Section 8. The appendix contains some algebraic details used in the derivation of Algorithm 1.

2. Model description

The Gaussian heteroscedastic nonparametric regression model has generic form:

$$y_i \stackrel{\text{ind.}}{\sim} \text{N}(f(x_i), g(x_i)), \quad 1 \leq i \leq n, \quad (2)$$

where (x_i, y_i) is the i th predictor/response pair of a regression data-set and $\stackrel{\text{ind.}}{\sim}$ means ‘distributed independently as’. The functions f and g are smooth, but otherwise arbitrary, functions. We refer to f as the *mean function* and g as the *variance function*. We also use the term *standard deviation function* for \sqrt{g} .

Mixed model-based penalized spline fitting of (2) (e.g. Ruppert *et al.* 2003) involves modelling f and g according to:

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^{K_u} u_k z_k^u(x), \quad u_k \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2)$$

(3)

$$\text{and } g(x) = \exp\left(\gamma_0 + \gamma_1 x + \sum_{k=1}^{K_v} v_k z_k^v(x)\right), \quad v_k \stackrel{\text{ind.}}{\sim} N(0, \sigma_v^2).$$

The $\{z_k^u : 1 \leq k \leq K_u\}$ and $\{z_k^v : 1 \leq k \leq K_v\}$ are spline bases of sizes K_u and K_v , respectively. Our default for the z_k^u and z_k^v are suitably transformed cubic O'Sullivan splines, as described in section 4 of Wand & Ormerod (2008). If $K_u = K_v$ then the two bases are identical, but we leave open the possibility for different basis sizes for the mean and variance functions.

In this article, we take a Bayesian approach to fitting (2) and (3), and impose the following priors on the model parameters:

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

with hyperparameters $\sigma_\beta, \sigma_\gamma, A_u, A_v > 0$ to be specified by the user. The Half-Cauchy (A) density function is given by $p(x) = \{2/(\pi A)\}/\{1 + (x/A)^2\}$, $x > 0$. Assuming that the data have been pre-transformed to have zero mean and unit variance, our default setting of the hyperparameters throughout this article are:

$$\sigma_\beta = \sigma_\gamma = A_u = A_v = 10^5.$$

Result 5 of Wand *et al.* (2011) allows us to replace $\sigma_u \sim \text{Half-Cauchy}(A_u)$ by

$$\sigma_u^2 | a_u \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_u\right), \quad a_u \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_u^2\right), \quad (4)$$

where $x \sim \text{Inverse-Gamma}(A, B)$ means that x has an Inverse Gamma distribution with shape parameter $A > 0$ and rate parameter $B > 0$. The corresponding density function is $p(x) = B^A \Gamma(A)^{-1} x^{-A-1} \exp(-B/x)$, $x > 0$. Representation (4) is more amenable to variational approximate inference. A similar replacement is made for σ_v .

The full Bayesian hierarchical model corresponding to (2) is:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{u}, \mathbf{v} &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_u \mathbf{u}, \text{diag}\{\exp(\mathbf{X}\boldsymbol{\gamma} + \mathbf{Z}_v \mathbf{v})\}), \\ \mathbf{u} | \sigma_u^2 &\sim N(\mathbf{0}, \sigma_u^2 \mathbf{I}), \quad \mathbf{v} | \sigma_v^2 \sim N(\mathbf{0}, \sigma_v^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}), \quad \boldsymbol{\gamma} \sim N(\mathbf{0}, \sigma_\gamma^2 \mathbf{I}), \\ \sigma_u^2 | a_u &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_u\right), \quad a_u \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_u^2\right), \\ \sigma_v^2 | a_v &\sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/a_v\right), \quad a_v \sim \text{Inverse-Gamma}\left(\frac{1}{2}, 1/A_v^2\right). \end{aligned} \quad (5)$$

Here $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are 2×1 vectors of fixed effects containing (β_0, β_1) and (γ_0, γ_1) respectively, \mathbf{u} is the $K_u \times 1$ vector whose entries are u_1, \dots, u_{K_u} , \mathbf{v} is a $K_v \times 1$ vector defined similarly, and σ_u^2 and σ_v^2 are variance components corresponding to \mathbf{u} and \mathbf{v} respectively. The design matrix \mathbf{X} is the $n \times 2$ matrix consisting of a column of ones and a column containing the x_i 's, $1 \leq i \leq n$. Also, \mathbf{Z}_u is the $n \times K_u$ matrix with (i, k) entry equal to $z_k^u(x_i)$ for $1 \leq k \leq K_u$ and \mathbf{Z}_v is the $n \times K_v$ matrix with (i, k) entry equal to $z_k^v(x_i)$ for $1 \leq k \leq K_v$.

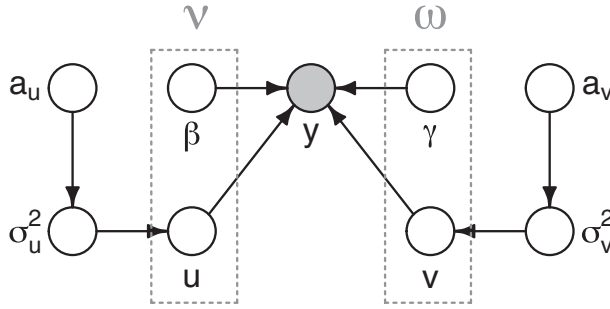


Figure 2. Directed acyclic graph for the model in (5). The shaded node corresponds to the observed data vector. Random effects and auxiliary variables are referred to as hidden nodes. The grey dashed boxes indicate that β and u are combined into a vector denoted by ν and ω is the concatenation of γ and v .

It is also convenient to combine the mean function coefficients and variance function coefficients into single vectors:

$$\nu \equiv \begin{bmatrix} \beta \\ u \end{bmatrix} \quad \text{and} \quad \omega \equiv \begin{bmatrix} \gamma \\ v \end{bmatrix}.$$

An analogous combining of the design matrices: $C_\nu \equiv [X Z_u]$, $C_\omega \equiv [X Z_v]$, then allows us to write

$$X\beta + Z_u u = C_\nu \nu \quad \text{and} \quad X\gamma + Z_v v = C_\omega \omega.$$

Figure 2 shows the directed acyclic graph corresponding to the model conveyed in (5).

MCMC schemes for fitting and inference in (5) are relatively straightforward to devise and implement. The BUGS and Stan MCMC-based inference engines also support (5) and illustration of this fact is given in Section 4. However, MCMC does not scale well to high volume/velocity data and larger models with heteroscedastic nonparametric regression components.

3. Variational inference methodology

We now consider the problem of mean field-type variational inference for (5) (e.g. Wainwright & Jordan 2008; Ormerod & Wand 2010). This involves an approximation to the joint posterior density function of the form

$$p(\nu, \omega, \sigma^2, \mathbf{a} | \mathbf{y}) \approx q(\nu) q(\omega) q(\sigma^2) q(\mathbf{a}) \tag{6}$$

where $\sigma^2 = (\sigma_u^2, \sigma_v^2)$ and $\mathbf{a} = (a_u, a_v)$. The q -densities are chosen to minimize the Kullback–Leibler divergence between the left-hand-side and right-hand-side of (6). This is equivalent to maximizing the variational lower bound on the marginal log-likelihood $\log p(\mathbf{y})$:

$$\log \underline{p}(\mathbf{y}; q) \equiv E_q [\log \{p(\nu, \omega, \sigma^2, \mathbf{a}, \mathbf{y})\} - \log \{q(\nu) q(\omega) q(\sigma^2) q(\mathbf{a})\}].$$

The optimal q densities, denoted by q^* , can be shown to satisfy (e.g. Ormerod & Wand 2010):

$$\begin{aligned}
 q^*(\mathbf{v}) &\propto \exp \{E_{q(-\mathbf{v})} \log p(\mathbf{v}|\text{rest})\}, & q^*(\boldsymbol{\omega}) &\propto \exp \{E_{q(-\boldsymbol{\omega})} \log p(\boldsymbol{\omega}|\text{rest})\}, \\
 q^*(\sigma^2) &\propto \exp \{E_{q(-\sigma^2)} \log p(\sigma^2|\text{rest})\} & \text{and} & q^*(\mathbf{a}) \propto \exp \{E_{q(-\mathbf{a})} \log p(\mathbf{a}|\text{rest})\},
 \end{aligned}
 \tag{7}$$

where, for example $E_{q(-\mathbf{v})}$ denotes expectation with respect to the q -densities of all parameters except \mathbf{v} . Also ‘rest’ denotes all of the random variables in the model other than those in \mathbf{v} , including \mathbf{y} .

The full conditionals $p(\mathbf{v}|\text{rest})$, $p(\sigma^2|\text{rest})$ and $p(\mathbf{a}|\text{rest})$ each have standard forms which result in closed form q -density expressions. However, $p(\boldsymbol{\omega}|\text{rest})$ is a non-standard form and challenging integrals arise in the determination of $q^*(\boldsymbol{\omega})$. We achieve a tractable solution by imposing the additional restriction that

$$q(\boldsymbol{\omega}) = q(\boldsymbol{\omega}; \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) \text{ is a } N(\boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) \text{ density function}
 \tag{8}$$

for some mean vector $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and covariance matrix $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$. The corresponding marginal log-likelihood is

$$\begin{aligned}
 \log \underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) &\equiv E_q \left[\log p(\mathbf{v}, \boldsymbol{\omega}, \sigma^2, \mathbf{a}, \mathbf{y}) \right. \\
 &\quad \left. - \log \{q(\mathbf{v})q(\boldsymbol{\omega}; \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})q(\sigma^2)q(\mathbf{a})\} \right]
 \end{aligned}
 \tag{9}$$

and minimization of the Kullback–Leibler divergence corresponds to $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ being chosen to maximise (9). This approach, due to Knowles & Minka (2011), has been labelled *non-conjugate variational message passing* since it provides a way of circumventing non-conjugacies in MFVB. (Note that *variational message passing* is an alternative formulation of MFVB, see e.g. Minka & Winn 2008.) Knowles & Minka (2011) propose a fixed point iteration scheme as a means of maximizing (9). Wand (2014) provides the algebraic details and simplification of fixed point updates for the Multivariate Normal q -density parameters, such as (8).

For fixed $(\boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})})$, application of (7), with $q^*(\boldsymbol{\omega})$ omitted, leads to:

$$\begin{aligned}
 q^*(\mathbf{v}) &\text{ is the } N(\boldsymbol{\mu}_{q(\mathbf{v})}, \boldsymbol{\Sigma}_{q(\mathbf{v})}) \text{ density function,} & q^*(\sigma^2) &= q^*(\sigma_u^2, \sigma_v^2) \text{ is a product} \\
 &\text{ of Inverse-Gamma} \left(\frac{1}{2}(K_u + 1), B_{q(\sigma_u^2)} \right) & & \text{and Inverse-Gamma} \left(\frac{1}{2}(K_v + 1), B_{q(\sigma_v^2)} \right) \\
 &\text{ density functions, and } q^*(\mathbf{a}) = q^*(a_u, a_v) & & \text{is a product of} \\
 &\text{ Inverse-Gamma}(1, B_{q(a_u)}) & & \text{and Inverse-Gamma}(1, B_{q(a_v)}) \text{ density functions}
 \end{aligned}
 \tag{10}$$

for parameters $\boldsymbol{\mu}_{q(\mathbf{v})}$ and $\boldsymbol{\Sigma}_{q(\mathbf{v})}$, the mean and covariance matrix of $q^*(\mathbf{v})$, $B_{q(\sigma_u^2)}$, the rate parameter of $q^*(\sigma_u^2)$, $B_{q(\sigma_v^2)}$, the rate parameter of $q^*(\sigma_v^2)$, $B_{q(a_u)}$, the rate parameter of $q^*(a_u)$ and $B_{q(a_v)}$, the rate parameter of $q^*(a_v)$. These optimal parameters are all inter-related and obtained through an iterative algorithm, listed below as Algorithm 1. The algorithm also includes fixed point iterative updates for $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$. Details on the derivation of (10), as well as the fixed point updates, are given in an Appendix.

Before presenting Algorithm 1 some additional notation is required. We use \leftarrow to denote assignment, as used in the R programming environment. For two matrices \mathbf{A} and \mathbf{B} of equal size, $\mathbf{A} \odot \mathbf{B}$ denotes their element-wise product. If \mathbf{a} is a $d \times 1$ vector then $\text{diag}(\mathbf{a})$

is the $d \times d$ diagonal matrix with the entries of \mathbf{a} along the diagonal. For a $d \times d$ matrix A , $\text{diagonal}(A)$ is the $d \times 1$ vector comprising the diagonal entries of A . Also, $\boldsymbol{\mu}_{q(\mathbf{u})}$ is defined to be the sub-vector of $\boldsymbol{\mu}_{q(\mathbf{v})}$ corresponding to \mathbf{u} , $\boldsymbol{\Sigma}_{q(\mathbf{u})}$ is the sub-matrix of $\boldsymbol{\Sigma}_{q(\mathbf{v})}$ corresponding to \mathbf{u} . The symbols $\boldsymbol{\mu}_{q(\mathbf{v})}$ and $\boldsymbol{\Sigma}_{q(\mathbf{v})}$ are defined similarly.

Algorithm 1. MFVB algorithm for the determination of the optimal parameters in $q^*(\boldsymbol{\omega})$, $q^*(\mathbf{v})$, $q^*(\sigma_u^2)$, $q^*(\sigma_v^2)$, $q^*(a_u)$ and $q^*(a_v)$.

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

Cycle:

$$\begin{aligned} \boldsymbol{\psi}_{q(\boldsymbol{\omega})} &\leftarrow \exp \left\{ -\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \frac{1}{2} \text{diagonal} \left(\mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \mathbf{C}_\omega^\top \right) \right\} \\ \boldsymbol{\Sigma}_{q(\mathbf{v})} &\leftarrow \left(\mathbf{C}_v^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \mathbf{C}_v + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right)^{-1} \\ \boldsymbol{\mu}_{q(\mathbf{v})} &\leftarrow \boldsymbol{\Sigma}_{q(\mathbf{v})} \mathbf{C}_v^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) \mathbf{y} \\ \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} &\leftarrow \left\{ \mathbf{C}_\omega^\top \text{diag} \left(\boldsymbol{\mu}_{q(r_v^2)} \odot \boldsymbol{\psi}_{q(\boldsymbol{\omega})} \right) \mathbf{C}_\omega + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \right\}^{-1} \\ \boldsymbol{\mu}_{q(\boldsymbol{\omega})} &\leftarrow \boldsymbol{\mu}_{q(\boldsymbol{\omega})} + \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \left\{ \mathbf{C}_\omega^\top \left(\boldsymbol{\mu}_{q(r_v^2)} \odot \boldsymbol{\psi}_{q(\boldsymbol{\omega})} \right) - \begin{bmatrix} \sigma_\gamma^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\mu}_{q(1/\sigma_v^2)} \mathbf{I} \end{bmatrix} \boldsymbol{\mu}_{q(\boldsymbol{\omega})} \right\} \\ \mu_{q(1/a_u)} &\leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^2) \quad ; \quad \mu_{q(1/a_v)} \leftarrow 1 / (\mu_{q(1/\sigma_v^2)} + A_v^2) \\ \mu_{q(1/\sigma_u^2)} &\leftarrow \frac{K_u + 1}{2\mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(\mathbf{u})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{u})})} \\ \mu_{q(1/\sigma_v^2)} &\leftarrow \frac{K_v + 1}{2\mu_{q(1/a_v)} + \|\boldsymbol{\mu}_{q(\mathbf{v})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\mathbf{v})})} \\ \boldsymbol{\mu}_{q(r_v^2)} &\leftarrow \text{diagonal} \left\{ (\mathbf{y} - \mathbf{C}_v \boldsymbol{\mu}_{q(\mathbf{v})}) (\mathbf{y} - \mathbf{C}_v \boldsymbol{\mu}_{q(\mathbf{v})})^\top + \mathbf{C}_v \boldsymbol{\Sigma}_{q(\mathbf{v})} \mathbf{C}_v^\top \right\} \end{aligned}$$

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

The symbol \leftarrow indicates ‘is assigned the value’ or ‘is replaced by’.

Convergence of Algorithm 1 can be monitored using the following explicit expression for the marginal log-likelihood lower bound:

$$\begin{aligned} \log p(\mathbf{y}; q, \boldsymbol{\mu}_{q(\boldsymbol{\omega})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}) &= \frac{1}{2} (K_u + K_v + 4) - \frac{n}{2} \log(2\pi) + \log \Gamma \left(\frac{1}{2} (K_u + 1) \right) \\ &\quad + \log \Gamma \left(\frac{1}{2} (K_v + 1) \right) - 2 \log(\pi) - \log(A_u) - \log(A_v) \\ &\quad - \frac{n}{2} \mathbf{1}^\top (\mathbf{C}_\omega \boldsymbol{\mu}_{q(\boldsymbol{\omega})}) - \frac{1}{2} (\mathbf{y} - \mathbf{C}_v \boldsymbol{\mu}_{q(\mathbf{v})})^\top \text{diag}(\boldsymbol{\psi}_{q(\boldsymbol{\omega})}) (\mathbf{y} - \mathbf{C}_v \boldsymbol{\mu}_{q(\mathbf{v})}) \\ &\quad - \frac{1}{2} \text{tr} \left(\mathbf{C}_v^\top \mathbf{C}_v \boldsymbol{\Sigma}_{q(\mathbf{v})} \right) - \log(\sigma_\beta^2) - \log(\sigma_\gamma^2) + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\mathbf{v})}| \\ &\quad + \frac{1}{2} \log |\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}| - \frac{1}{2\sigma_\beta^2} (\|\boldsymbol{\mu}_{q(\boldsymbol{\beta})}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\boldsymbol{\beta})})) \end{aligned}$$

$$\begin{aligned}
& -\frac{1}{2\sigma_\gamma^2} \left\{ \|\boldsymbol{\mu}_{q(\gamma)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(\gamma)}) \right\} - \frac{1}{2}(K_u + 1) \log(B_{q(\sigma_u^2)}) \\
& - \frac{1}{2}(K_v + 1) \log(B_{q(\sigma_v^2)}) - \log(\mu_{q(1/\sigma_u^2)} + A_u^{-2}) \\
& - \log(\mu_{q(1/\sigma_v^2)} + A_v^{-2}) + \mu_{q(1/\sigma_u^2)}\mu_{q(1/a_u)} + \mu_{q(1/\sigma_v^2)}\mu_{q(1/a_v)}.
\end{aligned}$$

Unlike ordinary MFVB, there is no guarantee that each iteration will lead to an increase in $\underline{p}(\mathbf{y}; q, \boldsymbol{\mu}_{q(\omega)}, \boldsymbol{\Sigma}_{q(\omega)})$. Thus the iterations should be stopped when the absolute relative change in its logarithm falls below a negligible amount.

In Figure 3 we return to the example of Figure 1, armed with our new MFVB methodology. The top panels show the fitted mean functions, according to (5), with pointwise 95% credible sets. The model was fitted using MFVB corresponding to Algorithm 1 and MCMC based on BUGS (Spiegelhalter *et al.* 2003) with a burn-in of 5000, kept sample of 5000 and thinning factor of 5. The much faster MFVB fit is seen to be in excellent agreement with its more computationally costly benchmark. The middle panels show the fitted standard deviation functions, corresponding to \sqrt{g} in the notation of (2), and corresponding 95% credible sets. The agreement between MFVB and MCMC is good, rather than excellent, for \sqrt{g} . The \sqrt{g} -standardized residual plots at the bottom of Figure 3 indicate proper accounting for the heteroscedasticity and agreement with the normality.

We also conducted a comprehensive check in order to confirm that the relative change in the approximate marginal log-likelihood does not lead to early stopping. Over several hundred runs, with data generated from different scenarios, we recorded Bayes estimates of f and g at the stopping point and again with iterations continuing 25% beyond the stopping point. The differences in the estimates were negligible.

4. Performance assessment

We conducted a comprehensive simulation study to assess the performance of Algorithm 1 in terms of inferential accuracy and computing time. Data were simulated according to model (2) with $n = 500$ and the x_i 's uniform on $(0, 1)$. The four mean and variance function pairs are listed in Table 1, where $\phi(\cdot; \mu, \sigma)$ and $\Phi(\cdot; \mu, \sigma)$ respectively denote the density and distribution functions of the Normal distribution with mean μ and standard deviation σ . We generated 100 data-sets for each of the functions shown in Table 1.

For each model corresponding to a new replication, inference was achieved using MFVB via Algorithm 1 and MCMC for comparison. The R programming environment was used for implementation of Algorithm 1. The MFVB iterations were terminated when the relative change in $\log \underline{p}(\mathbf{y}; q)$ fell below 10^{-7} . The MCMC procedure was performed using BUGS with sample sizes the same as for the example in Section 3. The following sections provide details on the accuracy of MFVB against the MCMC benchmark.

4.1. Assessment of accuracy

Algorithm 1 yields fast approximate inference for the model parameters, however it does not guarantee that an adequate level of accuracy will be achieved. Figure 4 provides an accuracy assessment of Algorithm 1 using side-by-side boxplots of the accuracy scores for the parameters of interest.

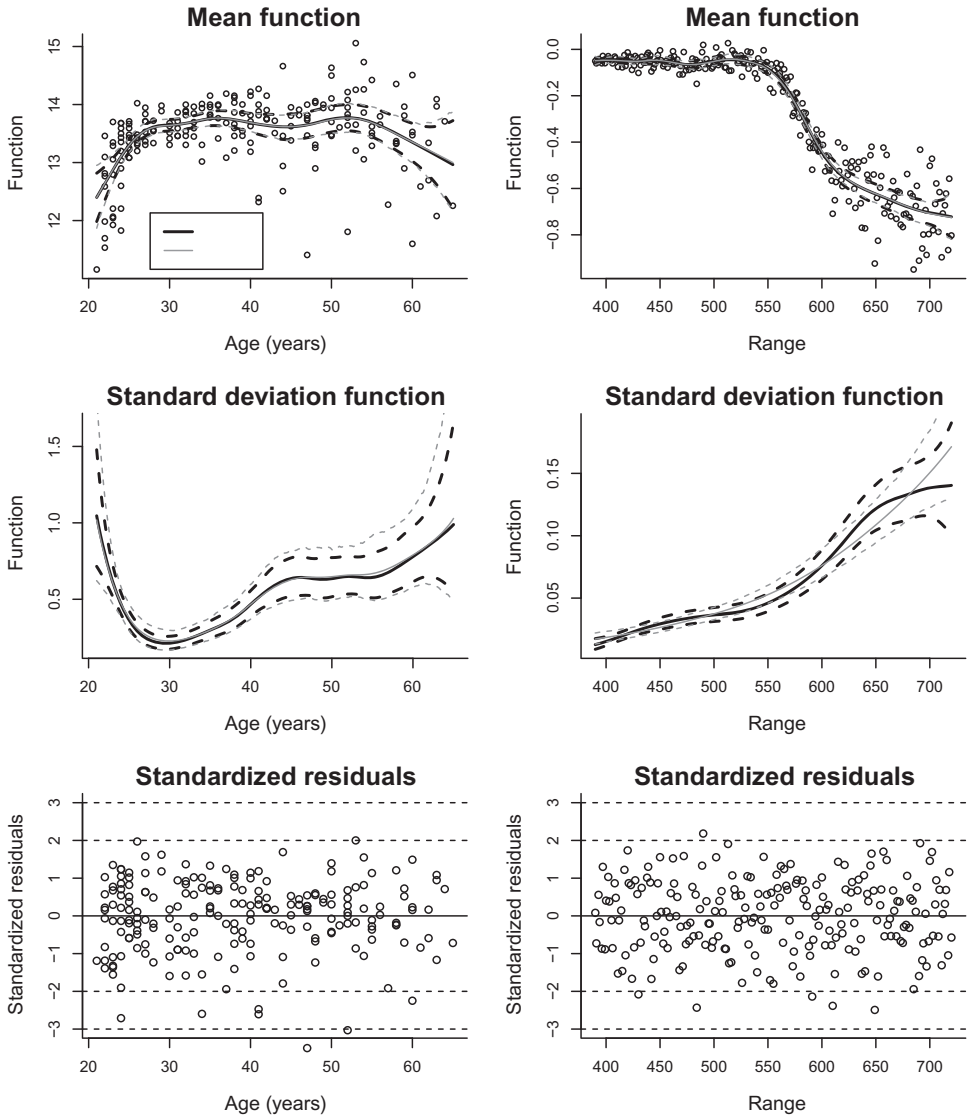


Figure 3. Top panels: fitted mean functions for two example data sets from Wand & Jones (1995) and Ruppert *et al.* (2003). The solid curves are approximate pointwise posterior means whilst the dashed curves are corresponding pointwise 95% credible sets. The approximate fits are based on MFVB via Algorithm 1 and MCMC via BUGS. Middle panels: similar to top panels but for the standard deviation function. Bottom panels: standardized residual plots based on $\{y - \hat{f}(x_i)\} / \sqrt{\{\hat{g}(x_i)\}}$ where \hat{f} and \hat{g} are the MFVB-approximate Bayes estimates of f and g .

For a generic parameter θ , the accuracy of $q^*(\theta)$ is defined to be

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

TABLE 1
Details of simulation study settings.

Setting	$f(x)$	$\log g(x)$
A	$\sin(3\pi x^2)$	$0.1 + \cos(4\pi x)$
B	$-1.02x + 0.018x^2 + 0.4\phi(x; 0.38, 0.08) + 0.08\phi(x; 0.75, 0.03)$	$-0.5 - \Phi(x; 0.2, 0.1) + 0.3x^2$
C	$0.35\phi(x; 0.01, 0.08) + 1.9\phi(x; 0.45, 0.23) + 1.8\{1 - \phi(x; 0.7, 0.14)\}$	$0.3\phi(x; 0, 0.2) + 0.4\phi(x; 1, 0.1)$
D	$\sin(3\pi x^2) - 1.02x + 0.018x^2 + 0.4\phi(x; 0.38, 0.08)$	$\cos(4\pi x) - 0.4 + 0.3x^2 - \Phi(x; 0.2, 0.1)$

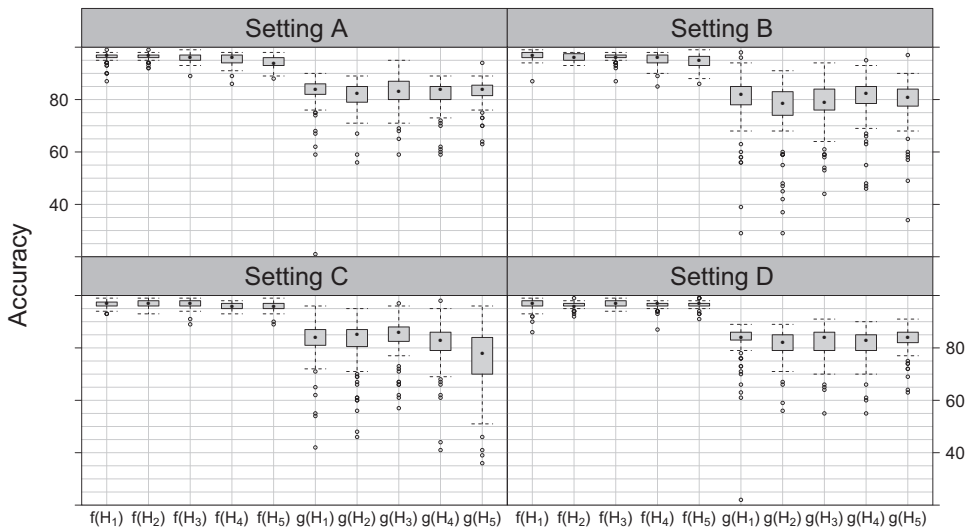


Figure 4. Summary of simulation study where the accuracy values are summarized as a boxplot.

Note that $p(\theta|y)$ can be approximated by using kernel density estimation applied to the MCMC sample. This accuracy score is justified in Faes *et al.* (2011). Accuracy is monitored for parameters $f(H_k)$ and $g(H_k)$, $1 \leq k \leq 5$, where the H_k are the sample hexiles of the x_i 's. The boxplots illustrate that most of the accuracies for the $f(H_k)$ lie around 90%, while accuracies for the $g(H_k)$ lie around 80%. These very good accuracy results are in keeping with the heuristics given in section 3.1 of Menictas & Wand (2013).

Figure 5 shows a comparison of MCMC and the MFVB fitted mean and standard deviation functions for the first replication in each of the four simulation settings. The MCMC and MFVB fits show excellent agreement for the mean functions and relatively good agreement for the standard deviation functions.

4.2. Assessment of coverage

We are also interested in the comparison between the coverage gained by the MFVB approximate credible intervals and the true coverage. Table 2 gives the percentages of the true parameter coverage based on the approximate 95% credible intervals attained from the

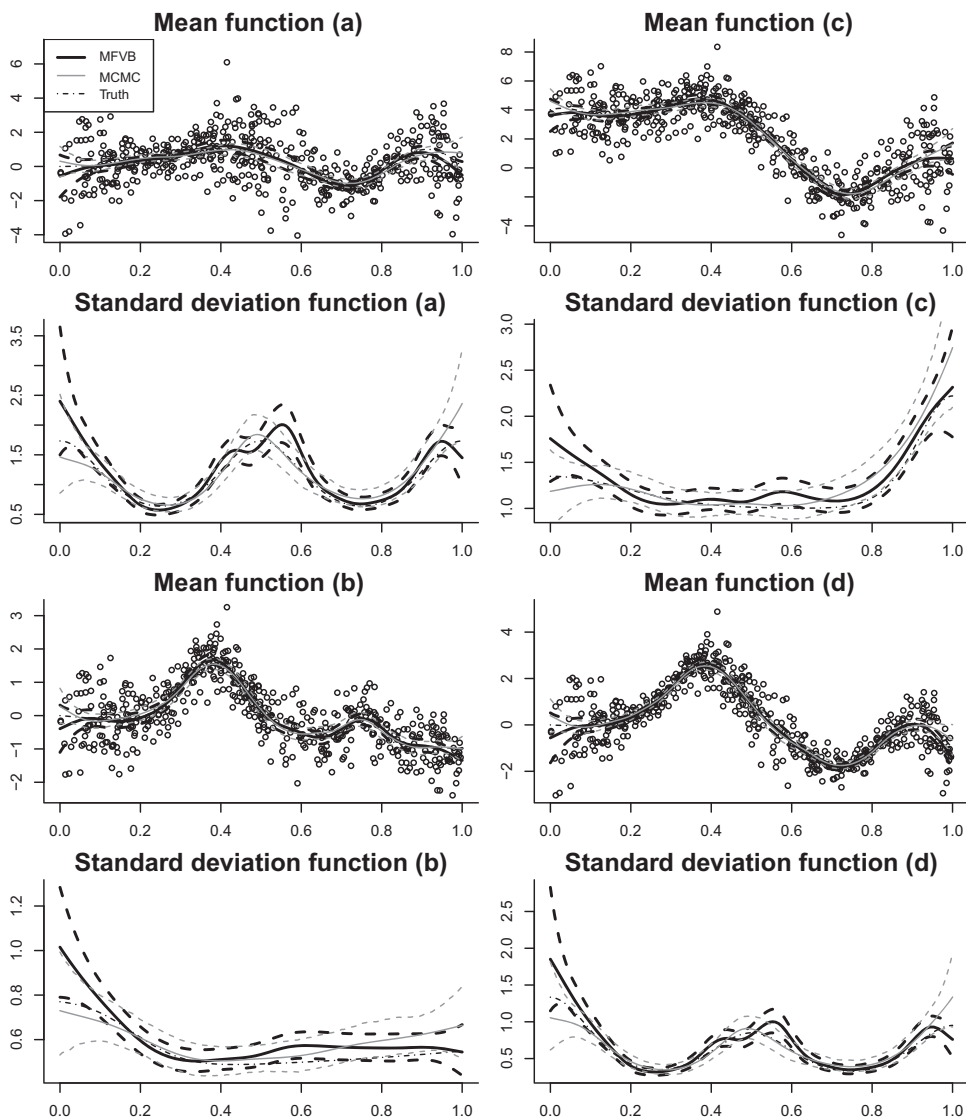


Figure 5. Comparison of MCMC and MFVB fitted functions for the first data-set in each of the four simulation settings.

MFVB posterior densities. The coverage overall is good and does not fall below 86%. As we have already seen in the previous section the performance of MFVB is excellent for the mean function and very good for the variance function.

4.3. Assessment of speed

During the running of the simulation we monitored the time taken per model to be fitted via MCMC and MFVB. The results are summarized in Table 3. The simulation was run on the first author’s desktop computer (Intel Core i5-2400 3.10 GHz processor, 8 GBytes of random access memory).

TABLE 2

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

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

TABLE 3

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

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

As mentioned previously, convergence was assessed differently for the two approaches. Also, the speed gains of MFVB are traded off against accuracy losses which are invoked by the product restriction given in (6). However, despite this concern, the results show that MFVB is approximately 200 times faster than MCMC when comparing all models. Thus we can assert that a model that takes minutes to run using MCMC, will take only seconds to run using our MFVB algorithm.

5. Extension to bivariate predictors

Algorithm 1 is relatively easy to extend to bivariate predictor nonparametric regression, which is closely related to *geostatistics* (e.g. Cressie 1993), which we briefly describe here. In this case, the data are of the form

$$(\mathbf{x}_i, y_i), \quad \mathbf{x}_i \in \mathbb{R}^2, y_i \in \mathbb{R}. \tag{11}$$

In the classical geostatistics scenario, the \mathbf{x}_i 's specify geographical locations. However, in (11) the \mathbf{x}_i 's could also represent pairs of non-geographical measurements.

The bivariate analogue of (2) is

$$y_i \overset{\text{ind.}}{\sim} N(f(\mathbf{x}_i), g(\mathbf{x}_i)), \quad 1 \leq i \leq n, \tag{12}$$

where f and g are real-valued functions on \mathbb{R}^2 . The extension of (3) is

$$f(\mathbf{x}) = \beta_0 + \beta_1^\top \mathbf{x} + \sum_{k=1}^{K_u} u_k z_k^u(\mathbf{x}), \quad u_k \overset{\text{ind.}}{\sim} N(0, \sigma_u^2) \tag{13}$$

and $g(\mathbf{x}) = \exp \left(\gamma_0 + \boldsymbol{\gamma}_1^\top \mathbf{x} + \sum_{k=1}^{K_v} v_k z_k^v(\mathbf{x}) \right), \quad v_k \overset{\text{ind.}}{\sim} N(0, \sigma_v^2),$

where each of β_1 and γ_1 are 2×1 vectors. The functions $\{z_k^u : 1 \leq k \leq K_u\}$ and $\{z_k^v : 1 \leq k \leq K_v\}$ are now bivariate spline basis functions. A reasonable default for the z_k^u (Ruppert *et al.* 2003) is the low-rank thin plate spline basis with k th element:

$$z_k^u(\mathbf{x}) = r(\|\mathbf{x} - \kappa_k^u\|) [r(\|\kappa_k^u - \kappa_{k'}^u\|)]^{-1/2}, \tag{14}$$

$1 \leq k, k' \leq K_u$

where $\kappa_1^u, \dots, \kappa_{K_u}^u$ is a set of bivariate knot locations that efficiently cover the space of the \mathbf{x}_i 's and $r(x) \equiv x^2 \log(x)$. The default z_k^v 's have an analogous definition.

According to this set-up, the only difference between the univariate nonparametric heteroscedastic regression model, treated in Section 3, and its bivariate counterpart is the basis functions and their coefficients. Hence, Algorithm 1 can be used to fit the bivariate nonparametric heteroscedastic regression model by replacing \mathbf{v} , ω , C_v and C_ω from Section 3 with

$$\mathbf{v} \equiv \begin{bmatrix} \beta_0 \\ \beta_1 \\ \mathbf{u} \end{bmatrix}, \quad \omega \equiv \begin{bmatrix} \gamma_0 \\ \gamma_1 \\ \mathbf{v} \end{bmatrix}, \quad C_v \equiv \left[\begin{array}{c} 1 \ \mathbf{x}_i^\top \\ z_k^u(\mathbf{x}_i) \end{array} \right]_{\substack{1 \leq k \leq K_u \\ 1 \leq i \leq n}} \quad \text{and} \quad C_\omega \equiv \left[\begin{array}{c} 1 \ \mathbf{x}_i^\top \\ z_k^v(\mathbf{x}_i) \end{array} \right]_{\substack{1 \leq k \leq K_v \\ 1 \leq i \leq n}}.$$

We fitted (13) to geo-referenced data on sea-floor sediment pollution in the North Sea (source: Pebesma & Duin 2005). The data are stored in the `pcb` data-frame within the R package `gstat` (Pebesma 2004). The response variable is a measurement of polychlorinated biphenyl with Ballschmitter-Zell congener number 138 (PCB-138). The motivating study is concerned with spatial and temporal variability of PCB-138. For the purposes of illustration, we ignore the temporal aspect and focus on geographical variability in the mean and variance of the response. In the notation of model (12)–(13), the variables are $\mathbf{x} = (x_1, x_2)$ where

- $x_1 = x$ -coordinate in the Universal Mobile Telecommunications System for Zone 31,
- $x_2 = y$ -coordinate in the Universal Mobile Telecommunications System for Zone 31,
- and
- $y = \text{PCB-138}$ measured on the sediment fraction smaller than 63 parts per million, in $\mu\text{g}/\text{kg}$ dry matter.

The sample size is $n = 216$ and $K_u = K_v = 50$ thin plate basis splines were used for each functional fit. The estimated mean and standard deviation functions are shown in Figure 6. Both functions are seen to exhibit pronounced spatial effects. Simple bivariate predictor models that ignore the heteroscedasticity described by the right panel of Figure 6 would lead to erroneous prediction intervals.

Higher dimensional heteroscedastic nonparametric regression can be achieved via Algorithm 1 with little notational change from the bivariate case treated here. The only required modification involves higher-dimensional thin plate spline basis functions instead of those given by (14).

6. Extension to additive models

The final semiparametric regression extension, that we discuss briefly in this section, is additive models with multiplicative variance functions. Models of this type have a small literature with Rigby & Stasinopoulos (2005) being a key reference. Their generic form is

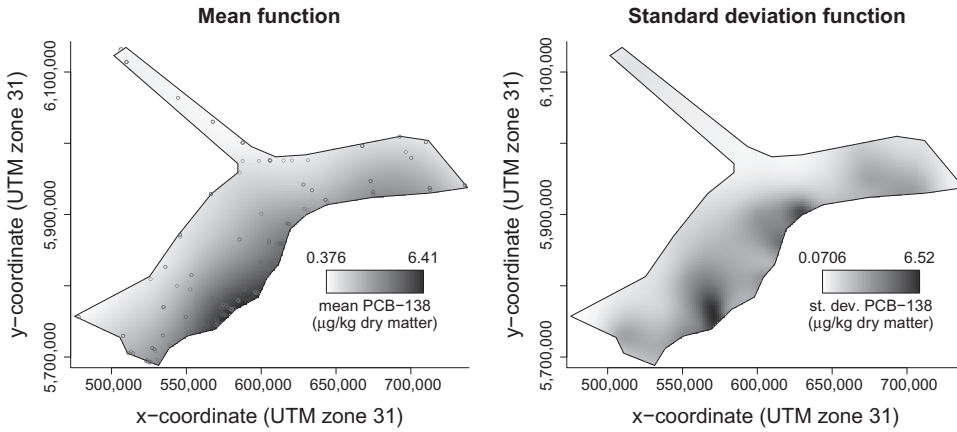


Figure 6. Left panel: Fitted mean function for the polychlorinated biphenyl data described in the text, based on MFVB via Algorithm 1. Right panel: similar to top panel but for the standard deviation function.

$$y_i \sim N\left(\beta_0 + \sum_{j=1}^d f_j(x_{ji}), \exp\left(\gamma_0 + \sum_{j=1}^d h_j(x_{ji})\right)\right), \quad 1 \leq i \leq n. \tag{15}$$

Here f_j and h_j , $1 \leq j \leq d$, are smooth but otherwise arbitrary functions. We use penalized spline models of the form

$$f_j(x) = \beta_j x + \sum_{k=1}^{K_j^u} u_{jk} z_{jk}^u(x), \quad u_{jk} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{uj}^2) \tag{16}$$

and $h_j(x) = \gamma_j x + \sum_{k=1}^{K_j^v} v_{jk} z_{jk}^v(x), \quad v_{jk} \stackrel{\text{ind.}}{\sim} N(0, \sigma_{vj}^2).$

The $\{z_{jk}^u : 1 \leq k \leq K_j^u\}$, $1 \leq j \leq d$, are spline bases of sizes K_j^u , analogous to those presented in Section 2. The $\{z_{jk}^v : 1 \leq k \leq K_j^v\}$, $1 \leq j \leq d$, are similarly defined. The priors on the regression coefficients and standard deviation parameters are

$$\begin{aligned} \beta_j &\stackrel{\text{ind.}}{\sim} N(0, \sigma_\beta^2), & \gamma_j &\stackrel{\text{ind.}}{\sim} N(0, \sigma_\gamma^2), \\ \sigma_{uj} &\stackrel{\text{ind.}}{\sim} \text{Half-Cauchy}(A_u), & \sigma_{vj} &\stackrel{\text{ind.}}{\sim} \text{Half-Cauchy}(A_v). \end{aligned} \tag{17}$$

The Bayesian model given by (15)–(17) admits a closed form non-conjugate MFVB algorithm, with the regression coefficients for the full mean and variance functions each being Multivariate Normal. Details and illustration are given in Menictas (2015). Section 5.2 of Wand (2014) also contains an illustration for data from a Californian air pollution study.

The extension to additive models with parametric components is trivial. For example, if binary predictor data b_i , $1 \leq i \leq n$, are also available then the following extension of (15):

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

can be handled via simple additions to the design matrices and coefficient vectors. Similar comments apply to the addition of continuous predictors that have purely linear impact on either the mean function or log-variance function.

7. Real-time heteroscedastic nonparametric regression

Almost all nonparametric regression methodology presented to date make the assumption that the data are processed in batch, that is, at the same time. However, some disadvantages of batch processing include the requirement that analysis must wait until the entire data set has been collected, and often the need to store the entire data set in memory. In the real-time case, the analysis is updated as each new data point is collected. This is beneficial, and sometimes essential, for both high volume and/or velocity data. In this section we present a variation of Algorithm 1 that allows for real-time heteroscedastic nonparametric regression.

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

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

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

$$\begin{aligned} \mathbf{C}_\nu &\leftarrow [\mathbf{C}_\nu^\top \mathbf{c}_{\nu, \text{new}}]^\top; \mathbf{C}_\omega \leftarrow [\mathbf{C}_\omega^\top \mathbf{c}_{\omega, \text{new}}]^\top; \mathbf{y} \leftarrow [\mathbf{y}^\top \mathbf{y}_{\text{new}}]^\top \\ \mu_{q(r_v^2)} &\leftarrow \text{diagonal} \left\{ (\mathbf{y} - \mathbf{C}_\nu \mu_{q(\nu)}) (\mathbf{y} - \mathbf{C}_\nu \mu_{q(\nu)})^\top + \mathbf{C}_\nu \Sigma_{q(\nu)} \mathbf{C}_\nu^\top \right\} \\ \psi_{q(\omega)} &\leftarrow \exp \left\{ -\mathbf{C}_\omega \mu_{q(\omega)} + \frac{1}{2} \text{diagonal} (\mathbf{C}_\omega \Sigma_{q(\omega)} \mathbf{C}_\omega^\top) \right\} \\ \Sigma_{q(\nu)} &\leftarrow \left(\mathbf{C}_\nu^\top \text{diag} (\psi_{q(\omega)}) \mathbf{C}_\nu + \begin{bmatrix} \sigma_\beta^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_u^2)} \mathbf{I}_{K_u} \end{bmatrix} \right)^{-1} \\ \mu_{q(\nu)} &\leftarrow \Sigma_{q(\nu)} \mathbf{C}_\nu^\top \text{diag} (\psi_{q(\omega)}) \mathbf{y} \\ \Sigma_{q(\omega)} &\leftarrow \left(\mathbf{C}_\omega^\top \text{diag} (\mu_{q(r_v^2)} \odot \psi_{q(\omega)}) \mathbf{C}_\omega + \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_v^2)} \mathbf{I}_{K_v} \end{bmatrix} \right)^{-1} \\ \mu_{q(\omega)} &\leftarrow \mu_{q(\omega)} + \Sigma_{q(\omega)} \left\{ \mathbf{C}_\omega^\top (\mu_{q(r_v^2)} \odot \psi_{q(\omega)} - \mathbf{1}) - \begin{bmatrix} \sigma_\gamma^{-2} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mu_{q(1/\sigma_v^2)} \mathbf{I}_{K_v} \end{bmatrix} \mu_{q(\omega)} \right\} \end{aligned}$$

$$\begin{aligned}\mu_{q(1/a_u)} &\leftarrow 1 / (\mu_{q(1/\sigma_u^2)} + A_u^{-2}); \mu_{q(1/a_v)} \leftarrow 1 / (\mu_{q(1/\sigma_v^2)} + A_v^{-2}) \\ \mu_{q(1/\sigma_u^2)} &\leftarrow (K_u + 1) / \{2\mu_{q(1/a_u)} + \|\boldsymbol{\mu}_{q(u)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(u)})\} \\ \mu_{q(1/\sigma_v^2)} &\leftarrow (K_v + 1) / \{2\mu_{q(1/a_v)} + \|\boldsymbol{\mu}_{q(v)}\|^2 + \text{tr}(\boldsymbol{\Sigma}_{q(v)})\}\end{aligned}$$

until the analysis is complete or data no longer available.

The symbol \leftarrow indicates ‘is assigned the value’ or ‘is replaced by’.

Algorithm 2 processes each new entry of \mathbf{y} , denoted by y_{new} , and its corresponding row of \mathbf{C}_v and \mathbf{C}_ω , denoted by $\mathbf{c}_{v,new}$ and $\mathbf{c}_{\omega,new}$, successively in real time. The starting values for the real-time procedure are determined by performing a sufficiently large batch fit. This is explained in more detail in section 2.1.1 of Luts *et al.* (2014).

The web-site `realtime-semiparametric-regression.net` features a movie that illustrates Algorithm 2 for data simulated according to setting D . The warm-up sample size is $n_{warm} = 500$. The link for the movie is titled `Heteroscedastic nonparametric regression` and portrays the effectiveness of real time processing for mean and standard deviation function fitting.

8. Concluding remarks

We have developed closed form algorithms for fast batch and real-time fitting and inference for a variety of heteroscedastic semiparametric regression models.

The methodology also applies to larger models courtesy of the locality property of mean field variational inference methods. The new methodology has been shown to perform very well on simulated and actual data.

Appendix A: Derivation of optimal q -density functions

A.1. Derivation of $q^*(\mathbf{v})$

First note that

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

where ‘const’ denotes terms not depending on the argument of q^* . The form of $q^*(\mathbf{v})$ follows from this and the fact that

$$E_q(e^{-\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)}}) = \exp\{-\mathbf{C}_\omega \boldsymbol{\mu}_{q(\omega)} + \frac{1}{2} \mathbf{C}_\omega \boldsymbol{\Sigma}_{q(\omega)} \mathbf{C}_\omega^\top\}.$$

Therefore

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

and

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

A.2. Derivation of $q^*(\sigma_u^2)$ and $q^*(\sigma_v^2)$

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

The form of $q^*(\sigma_u^2)$ follows from the fact that

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

The expression for $\mu_{q(1/\sigma_u^2)}$ follows from a suitable result for the Inverse-Gamma distribution. For example, if random variable v has an Inverse-Gamma distribution with density function $p(v) = B^A \Gamma(A)^{-1} v^{-A-1} \exp(-v/B)$, then $E(1/v) = A/B$. Therefore

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

and

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

The derivation of $B_{q(\sigma_v^2)}$ and $\mu_{q(1/\sigma_v^2)}$ is similar.

A.3. Derivation of $q^*(a_u)$ and $q^*(a_v)$

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

The expressions for $B_{q(a_u)}$ and $\mu_{q(1/a_u)}$ follow immediately. The derivation of $B_{q(a_v)}$, and $\mu_{q(1/a_v)}$ is similar to that just shown above.

A.4. Derivation of the $\boldsymbol{\mu}_{q(\boldsymbol{\omega})}$ and $\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})}$ updates

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

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

where

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

D denotes derivative vector, as defined in Magnus & Neudecker (1999), and vec and vec^{-1} are as defined in Wand (2014). However, a result in the appendix of Opper & Archambeau (2009) shows that an equivalent expression for the first of these updates is

$$\boldsymbol{\Sigma}_{q(\boldsymbol{\omega})} \leftarrow (-\text{H}_{\boldsymbol{\mu}_{q(\boldsymbol{\omega})}} S)^{-1}$$

where \mathbf{H} denotes the Hessian matrix as defined in Magnus & Neudecker (1999). We work with this alternative form here. An explicit expression for S is

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

This leads to

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

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

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

Next,

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

Then, using $\text{diag}(a)b = a \odot b$,

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

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

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

Combining these expressions leads to the updates in Algorithm 1.

References

- BARBER, D. & BISHOP, C.M. (1997). Ensemble learning for multi-layer networks. In *Advances in Neural Information Processing Systems*, vol. 10, eds. M.I. JORDAN, K.J. KEARNS and S.A. SOLLA, pp. 395–401. Cambridge, MA: MIT Press.
- CHALLIS, E. & BARBER, D. (2013). Gaussian Kullback–Leibler approximate inference. *J. Mach. Learn. Res.* **14**, 2239–2286.
- CRAINICEANU, C.M., RUPPERT, D., CARROLL, R.J., JOSHI, A. & GOODNER, B. (2007). Spatially adaptive Bayesian penalized splines with heteroscedastic errors. *J. Comp. Graph. Stat.* **16**, 265–288.
- CRESSIE N.A.C. (1993). *Statistics for Spatial Data*, Revised Edition. New York: Wiley.
- FAES, C., ORMEROD, J.T. & WAND, M.P. (2011). Variational Bayesian inference for parametric and non-parametric regression with missing data. *J. Amer. Statist. Assoc.* **106**, 959–971.
- HINTON, G.E. & VAN CAMP, D. (1993). Keeping neural networks simple by minimizing the description length of the weights. In *Proceedings of the 6th Annual Association Computing Machinery Conference Computational Learning Theory*, pp. 5–13, New York: The Association for Computing Machinery.
- KNOWLES, D.A & MINKA, T.P. (2011). Non-conjugate variational message passing for multinomial and binary regression. In *Advances in Neural Information Processing Systems*, vol. 24, eds. J. SHAWE-TAYLOR, R.S. ZAMEL, P. BARTLETT, F. PEREIRA and K.Q. WEINBERGER, pp. 1701–1709. Cambridge, MA: MIT Press.
- LÁZARO-GREDILLA, M. & TITSIAS, M.K. (2011). Variational heteroscedastic Gaussian process regression. *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, Washington, USA.
- LUTS, J., BRODERICK, T. & WAND, M.P. (2014). Real-time semiparametric regression. *J. Comp. Graph. Stat.* **23**, 589–615.
- LUTS, J., WANG, S.S.J., ORMEROD, J.T. & WAND, M.P. (2015). Semiparametric regression analysis via Infer.NET. Unpublished manuscript.
- MAGNUS, J.R. & NEUDECKER, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Revised Edition. Chichester UK: Wiley.
- MARLEY, J.K. & WAND, M.P. (2010). Non-standard semiparametric regression via BRugs. *J. Stat. Softw.* **37**, 1–30.
- MENICTAS, M. (2015). Variational inference for heteroscedastic and longitudinal regression models (PhD Thesis). University of Technology, Sydney.
- MENICTAS, M. & WAND, M.P. (2013). Variational inference for marginal longitudinal semiparametric regression. *Statistics* **2**, 61–71.

- MINKA, T. & WINN, J. (2008). Gates: A graphical notation for mixture models. Microsoft Research Technical Report Series, MSR-TR-2008-185, 1–16.
- MINKA, T., WINN, J., GUIVER, J. & KNOWLES, D. (2014). Infer.NET 2.6. Available from URL: <http://research.microsoft.com/infernet> [Last accessed 26 Feb 2015.]
- NOTT, D.J., TRAN, M.-N. & LENG, C. (2012). Variational approximation for heteroscedastic linear models and matching pursuit algorithms. *Stat. Comput.* **22**, 497–512.
- OPPER, M. & ARCHAMBEAU, C. (2009). The variational Gaussian approximation revisited. *Neural Comput.* **21**, 786–792.
- ORMEROD, J.T. & WAND, M.P. (2010). Explaining variational approximations. *Amer. Statist.* **64**, 140–153.
- PEBESMA, E.J. (2004). Multivariate geostatistics in R: The *gstat* package. *Comput. Geosci.* **30**, 683–691.
- PEBESMA, E.J. & DUIN, R.N.M. (2005). Spatial patterns of temporal change in North Sea sediment quality on different spatial scales. In *Geostatistics for Environmental Applications: Proceedings of the Fifth European Conference on Geostatistics for Environmental Applications*, eds. P. RENARD, H. DEMOUGEOUT-RENARD and R. FROIDEVAUX, pp. 367–378. Berlin, Heidelberg: Springer-Verlag.
- R Core Team (2015). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, ISBN: 3-900051-07-0. Available from URL: <http://www.R-project.org> [Last accessed 26 Feb 2015.]
- RAIKO, T., VALPOLA, H., HARVA, M. & KARHUNEN, J. (2007). Building blocks for variational Bayesian learning of latent variable models. *J. Mach. Learn. Res.* **8**, 155–201.
- RIGBY, R.A. & STASINOPOULOS, D.M. (2005). Generalized additive models for location, scale and shape. *J. R. Statist. Soc. Ser. C. Appl. Statist.* **54**, 507–554.
- RUPPERT, D., WAND, M.P. & CARROLL, R.J. (2003). *Semiparametric Regression*. New York: Cambridge University Press.
- RUPPERT, D., WAND, M.P. & CARROLL, R.J. (2009). Semiparametric regression during 2003–2007. *Electron. J. Stat.* **3**, 1193–1256.
- SPIEGELHALTER, D.J., THOMAS, A., BEST, N.G., GILKS, W.R. & LUNN, D. (2003). BUGS: Bayesian Inference Using Gibbs Sampling. Medical Research Council Biostatistics Unit, Cambridge, UK. Available from URL: <http://www.mrc-bsu.cam.ac.uk/bugs> [Last accessed 26 Feb 2015.]
- Stan Development Team. (2013). Stan: A C++ library for probability and sampling. Version 1.3. Available from URL: <http://mc-stan.org/> [Last accessed 26 Feb 2015.]
- WAINWRIGHT, M.J. & JORDAN, M.I. (2008). Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.* **1**, 1–305.
- WAND, M.P. (2009). Semiparametric regression and graphical models. *Aust. N. Z. J. Stat.* **51**, 9–41.
- WAND, M.P. (2014). Fully simplified Multivariate Normal updates in non-conjugate variational message passing. *J. Mach. Learn. Res.* **15**, 1351–1369.
- WAND, M.P. & JONES, M.C. (1995). *Kernel Smoothing*. London: Chapman and Hall.
- WAND, M.P. & ORMEROD, J.T. (2008). On semiparametric regression with O’Sullivan penalised splines. *Aust. N. Z. J. Stat.* **50**, 179–198.
- WAND, M.P., ORMEROD, J.T., PADOAN, S.A. & FRÜHWIRTH, R. (2011). Mean field variational Bayes for elaborate distributions. *Bayesian Anal.* **6**, 847–900.