



Journal of Statistical Software

October 2018, Volume 87, Issue 2.

doi: [10.18637/jss.v087.i02](https://doi.org/10.18637/jss.v087.i02)

Semiparametric Regression Analysis via Infer.NET

Jan Luts TheSearchParty.com	Shen S. J. Wang University of Queensland	John T. Ormerod University of Sydney	Matt P. Wand University of Technology Sydney
---------------------------------------	---	---	---

Abstract

We provide several examples of Bayesian semiparametric regression analysis via the Infer.NET package for approximate deterministic inference in Bayesian models. The examples are chosen to encompass a wide range of semiparametric regression situations. Infer.NET is shown to produce accurate inference in comparison with Markov chain Monte Carlo via the BUGS package, but to be considerably faster. Potentially, this contribution represents the start of a new era for semiparametric regression, where large and complex analyses are performed via fast Bayesian inference methodology and software, mainly being developed within Machine Learning.

Keywords: additive mixed models, expectation propagation, generalized additive models, measurement error models, mean field variational Bayes, missing data models, penalized splines, variational message passing.

1. Introduction

Infer.NET (Minka, Winn, Guiver, and Knowles 2014) is a relatively new software package for performing approximate inference in large Bayesian models, using fast deterministic algorithms such as expectation propagation (Minka 2001) and variational message passing (Winn and Bishop 2005). We demonstrate its application to Bayesian semiparametric regression. A variety of situations are covered: non-Gaussian response, longitudinal data, bivariate functional effects, robustness, sparse signal penalties, missingness and measurement error.

The Infer.NET project is still in its early years and, and at the time of this writing, has not progressed beyond beta releases. We anticipate continual updating and enhancement for many years to come. In this article we are necessarily restricted to the capabilities of Infer.NET at the time of preparation. All of our examples use Infer.NET 2.7, which was released in March 2018.

The graphical model representation viewpoint of semiparametric regression (Wand 2009), and other advanced statistical techniques, has the attraction that various embellishments of the standard models can be accommodated by enlarging the underlying directed acyclic graph. For example, nonparametric regression with a missing predictor data model involves the addition of nodes and edges to the graph, corresponding to the missing data mechanism. Figure 4 of Faes, Ormerod, and Wand (2011) provides some specific illustrations. Marley and Wand (2010) exploited the graphical model representation viewpoint of semiparametric regression to handle a wide range of non-standard models via the Markov chain Monte Carlo inference engine implemented by BUGS (Lunn, Thomas, Best, and Spiegelhalter 2000). However, many of their examples take between several minutes and hours to run. The inference engines provided by *Infer.NET* allow much faster fitting for many common semiparametric regression models. On the other hand, BUGS is the more versatile package and not all models that are treated in Marley and Wand (2010) are supported by *Infer.NET*. Like BUGS the package *Infer.NET* allows the shifts from model estimation to model appropriateness by automatic deviation of fast approximate inference algorithms.

Semiparametric regression, summarized by Ruppert, Wand, and Carroll (2003, 2009), is a large branch of Statistics that includes nonparametric regression, generalized additive models, generalized additive mixed models, curve-by-factor interaction models, wavelet regression and geosadditive models. Parametric models such as generalized linear mixed models are special cases of semiparametric regression. Antecedent research for special cases such as nonparametric regression was conducted in the second half of the 20th Century at a time when data sets were smaller, computing power was lower and the Internet was either non-existent or in its infancy. Now, as we are in the mid-2010s, semiparametric regression is continually being challenged by the size, complexity and, in some applications, arrival speed of data sets requiring analysis. Implementation and computational speed are major limiting factors. This contribution represents the potential for a new era for semiparametric regression – tapping into 21st Century Machine Learning research on fast approximate inference on large Bayesian models (e.g., Minka 2001; Winn and Bishop 2005; Minka 2005; Minka and Winn 2008; Knowles and Minka 2011) and ensuing software development.

Section 2 lays down the definitions and notation needed to describe the models given in later sections and describes a set of support files for compiling and running *Infer.NET* programs. Each of Sections 3–10 illustrates a different type of semiparametric regression analysis via *Infer.NET*. All code is available in the supplementary material to this article. For most of the examples, we also compare the *Infer.NET* results with those produced by BUGS. Section 11 compares BUGS with *Infer.NET* in terms of versatility and computational speed. A summary of our findings on semiparametric analysis via *Infer.NET* is given in Section 12. Appendix A gives a detailed description of using *Infer.NET* to fit the simple semiparametric regression model of Section 3.

2. Preparatory infrastructure

The semiparametric regression examples rely on mathematical infrastructure and notation, which we describe in this section.

Distribution	Density function in x	Abbreviation
Normal	$(2\pi\sigma^2)^{-1/2} \exp\{-(x - \mu)^2/(2\sigma^2)\}; \quad \sigma > 0$	$N(\mu, \sigma^2)$
Laplace	$(2\sigma)^{-1} \exp(- x - \mu /\sigma); \quad \sigma > 0$	Laplace(μ, σ^2)
t	$\frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\pi\nu\sigma^2}\Gamma(\nu/2)\left\{1 + \frac{(x-\mu)^2}{\nu\sigma^2}\right\}^{\frac{\nu+1}{2}}}; \quad \sigma, \nu > 0$	$t(\mu, \sigma^2, \nu)$
Gamma	$\frac{B^A x^{A-1} e^{-Bx}}{\Gamma(A)}; \quad x > 0; A, B > 0$	Gamma(A, B)
Half-Cauchy	$\frac{2\sigma}{\pi(x^2 + \sigma^2)}; \quad x > 0; \sigma > 0$	Half-Cauchy(σ)

Table 1: Distributions used in the examples. The density function argument x and parameters range over \mathbb{R} unless otherwise specified.

2.1. Distributional notation

The density function of a random vector x is denoted by $p(x)$. The notation for a set of independent random variables y_i ($1 \leq i \leq n$) with distribution D_i is $y_i \stackrel{\text{ind.}}{\sim} D_i$. Table 1 lists the distributions used in the examples and the parametrization of their density functions.

2.2. Standardization and default priors

In real data examples, the measurements are often recorded on several different scales. Therefore, all continuous variables are standardized prior to analysis using Infer.NET or Markov chain Monte Carlo (MCMC) via BUGS. This transformation makes the analyses scale-invariant and can also lead to better behavior of the MCMC sampling.

Since we do not have prior knowledge about the model parameters in each of the examples, non-informative priors are used. The prior distributions for a fixed effects parameter vector β and a standard deviation parameter σ are

$$\beta \sim N(0, \tau_\beta^{-1}I), \quad \sigma \sim \text{Half-Cauchy}(A)$$

with default hyperparameters

$$\tau_\beta = 10^{-10}, \quad A = 10^5.$$

This is consistent with the recommendations given in Gelman (2006) for achieving non-informativeness for variance parameters. Infer.NET and BUGS do not offer direct specification of half-Cauchy distributions and therefore we use the result:

$$\begin{aligned} \text{If } x|a &\sim \text{Gamma}\left(\frac{1}{2}, a\right) \quad \text{and} \quad a \sim \text{Gamma}\left(\frac{1}{2}, 1/A^2\right), \\ \text{then } x^{-1/2} &\sim \text{Half-Cauchy}(A). \end{aligned} \tag{1}$$

This result allows for the imposition of a half-Cauchy prior using only gamma distribution specifications. Both *Infer.NET* and BUGS support the gamma distribution. This representation of the half-Cauchy prior suggests that the choice of A has little effect on the results provided the chosen value of A is sufficiently large which corresponds to the case where we have a diffuse prior on σ .

2.3. Variational message passing

Infer.NET has two inference engines, variational message passing (Winn and Bishop 2005) and expectation propagation (Minka 2001), for performing fast deterministic approximate inference in Bayesian models. Succinct summaries of variational message passing and expectation propagation are provided in Appendices A and B of Minka and Winn (2008).

Generally speaking, variational message passing is more amenable to semiparametric regression than expectation propagation. It is a special case of mean field variational Bayes (e.g., Wainwright and Jordan 2008). The essential idea of mean field variational Bayes is to approximate joint posterior density functions such as $p(\theta_1, \theta_2, \theta_3 | D)$, where D denotes the observed data, by product density forms such as

$$q_{\theta_1}(\theta_1) q_{\theta_2}(\theta_2) q_{\theta_3}(\theta_3), \quad q_{\theta_1, \theta_3}(\theta_1, \theta_3) q_{\theta_2}(\theta_2), \quad \text{or} \quad q_{\theta_1}(\theta_1) q_{\theta_2, \theta_3}(\theta_2, \theta_3). \quad (2)$$

The choice of the product density form is usually made by trading off tractability against minimal imposition of product structure. Once this choice is made, the optimal q -density functions are chosen to minimize the Kullback-Leibler divergence from the exact joint posterior density function. For example, if the second product form in (2) is chosen then the optimal density functions $q_{\theta_1, \theta_3}^*(\theta_1, \theta_3)$ and $q_{\theta_2}^*(\theta_2)$ are those that minimize

$$\iiint q_{\theta_1, \theta_3}(\theta_1, \theta_3) q_{\theta_2}(\theta_2) \log \left\{ \frac{p(\theta_1, \theta_2, \theta_3 | D)}{q_{\theta_1, \theta_3}(\theta_1, \theta_3) q_{\theta_2}(\theta_2)} \right\} d\theta_1 d\theta_2 d\theta_3, \quad (3)$$

where the integrals range over the parameter spaces of θ_1, θ_2 and θ_3 . This minimization problem gives rise to an iterative scheme which, typically, has closed form updates and good convergence properties. A by-product of the iterations is a lower-bound approximation to the marginal likelihood $p(D)$, which we denote by $\underline{p}(D; q)$. The iterative scheme of q -density updates can be shown to monotonically increase the lower bound $\underline{p}(D; q)$. These updates are guaranteed to converge to at least a local maximizer of $\underline{p}(D; q)$. Issues associated with multiple local maximizers of $\underline{p}(D; q)$ typically arise when performing inference using discrete random variables as model components. None of the models considered in this paper fall into this category.

Further details on, and several examples of, mean field variational Bayes are provided by Section 2.2 of Ormerod and Wand (2010). Each of these examples can also be expressed, equivalently, in terms of variational message passing.

In typical semiparametric regression models the subscripting on the q -density functions is quite cumbersome. Hence, it is suppressed for the remainder of the article. For example, $q(\theta_1, \theta_3)$ is taken to mean $q_{\theta_1, \theta_3}(\theta_1, \theta_3)$.

2.4. Expectation propagation

Expectation propagation is also based on product density restrictions such as (2), but differs in its method of obtaining the optimal density functions. It works with the reverse version of the Kullback-Leibler divergence than that given in (3), i.e.,

$$\iiint p(\theta_1, \theta_2, \theta_3 | D) \log \left\{ \frac{q_{\theta_1, \theta_3}(\theta_1, \theta_3) q_{\theta_2}(\theta_2)}{p(\theta_1, \theta_2, \theta_3 | D)} \right\} d\theta_1 d\theta_2 d\theta_3, \quad (4)$$

and, therefore, leads to different iterative algorithms and approximating density functions.

Expectation propagation overcomes the problem of minimizing (4) via a Kullback-Leibler projection onto exponential density functions. Minka (2005) develops a strategy for approximate minimization of (4) for general $p(\theta|D)$ and a given product restriction for $q(\theta)$ in terms of messages passed on an appropriate factor graph. A detailed explanation of expectation propagation accessible to statisticians is given in Kim and Wand (2016).

2.5. Mixed model-based penalized splines

Mixed model-based penalized splines are a convenient way to model nonparametric functional relationships in semiparametric regression models, and are amenable to the hierarchical Bayesian structures supported by Infer.NET. The penalized spline of a regression function f , with mixed model representation, takes the generic form

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k z_k(x), \quad u_k \stackrel{\text{ind.}}{\sim} N(0, \tau_u^{-1}). \quad (5)$$

Here $z_1(\cdot), \dots, z_K(\cdot)$ is a set of spline basis functions and τ_u controls the amount of penalization of the spline coefficients u_1, \dots, u_k . Throughout this article, we use O'Sullivan splines for the $z_k(\cdot)$. Wand and Ormerod (2008) provide details on their construction. O'Sullivan splines lead to (5) being a low-rank version of smoothing splines, which is also used by the R (R Core Team 2018) function `smooth.spline()`. The underlying O'Sullivan spline basis can be represented as a linear combination of B -splines with a fixed number of knots. Wand and Ormerod (2008) recommend quantile spaced knots, while Eilers and Marx (1996) use equally spaced knots. This choice has little effect on the fitted function provided a sufficiently large number of knots are used. A reasonable rule of thumb for the number of knots K is $K = \min(n_U/4, 35)$ where n_U is the number of unique x_i s. Univariate thin plate splines (Wood 2003) (multivariate thin plate splines are discussed below) and other types of radial basis functions (Fasshauer 2007) are alternative choices of basis functions that do not involve B -splines which could also be used here.

The simplest semiparametric regression model is the Gaussian response nonparametric regression model

$$y_i \stackrel{\text{ind.}}{\sim} N(f(x_i), \sigma_\varepsilon^2), \quad (6)$$

where (x_i, y_i) , $1 \leq i \leq n$ are pairs of measurements on continuous predictor and response variables. Mixed model-based penalized splines give rise to hierarchical Bayesian models for

(6) such as

$$\begin{aligned} y_i | \beta_0, \beta_1, u_1, \dots, u_K, \tau_\varepsilon &\stackrel{\text{ind.}}{\sim} N\left(\beta_0 + \beta_1 x_i + \sum_{k=1}^K u_k z_k(x_i), \tau_\varepsilon^{-1}\right), \\ u_k | \tau_u &\stackrel{\text{ind.}}{\sim} N(0, \tau_u^{-1}), \quad \beta_0, \beta_1 \stackrel{\text{ind.}}{\sim} N(0, \tau_\beta^{-1}), \\ \tau_u^{-1/2} &\sim \text{Half-Cauchy}(A_u), \quad \tau_\varepsilon^{-1/2} \sim \text{Half-Cauchy}(A_\varepsilon), \end{aligned} \quad (7)$$

where $\sigma_\varepsilon^2 = \tau_\varepsilon^{-1}$.

Such models allow nonparametric regression to be performed using Bayesian inference engines such as BUGS and *Infer.NET*. Our decision to work with precision parameters, rather than the variance parameters (which are common in the semiparametric regression literature), is driven by the former being the standard parametrization in BUGS and *Infer.NET*.

It is convenient to express (7) using matrix notation. This entails putting

$$\mathbf{y} \equiv \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} \mathbf{1} & \mathbf{x}_1 \\ \vdots & \vdots \\ \mathbf{1} & \mathbf{x}_n \end{bmatrix}, \quad \mathbf{Z} \equiv \begin{bmatrix} z_1(\mathbf{x}_1) & \cdots & z_K(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ z_1(\mathbf{x}_n) & \cdots & z_K(\mathbf{x}_n) \end{bmatrix} \quad (8)$$

and

$$\boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \text{and} \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_K \end{bmatrix}. \quad (9)$$

We then re-write (7) as

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon &\stackrel{\text{ind.}}{\sim} N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}), \\ \mathbf{u} | \tau_u &\sim N(\mathbf{0}, \tau_u^{-1}\mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), \\ \tau_u^{-1/2} &\sim \text{Half-Cauchy}(A_u), \quad \tau_\varepsilon^{-1/2} \sim \text{Half-Cauchy}(A_\varepsilon). \end{aligned}$$

The effective degrees of freedom (edf) corresponding to (7) are defined to be given by

$$\text{edf}(\tau_u, \tau_\varepsilon) \equiv \text{tr} \left([C^\top C + \text{blockdiag}\{\tau_\beta^2 \mathbf{I}_2, (\tau_u/\tau_\varepsilon) \mathbf{I}\}]^{-1} C^\top C \right) \quad (10)$$

and are a scale-free measure of the amount of fitting being performed by the penalized splines. Further details on effective degrees of freedom for penalized splines are given in Section 3.13 of [Ruppert *et al.* \(2003\)](#).

A directed acyclic graph representation of this model is given in [Figure 1](#). Random variables or vectors correspond to nodes while directed edges (i.e., arrows) convey conditional dependence. Nodes filled gray correspond to observed model variables (sometimes called evidence nodes), while unfilled nodes correspond to unobserved model variables (sometimes called hidden nodes). A further discussion of these graphical model representations can be found in [Wand \(2009\)](#) or [Faes *et al.* \(2011\)](#).

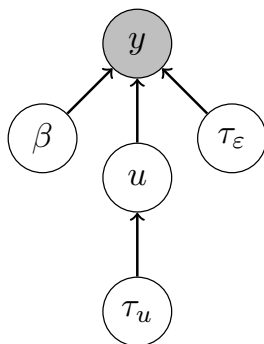


Figure 1: Directed acyclic graph representation of the mixed model-based penalized splines described in Section 2.5. Shading indicates observed data.

Extension to bivariate predictors

The bivariate predictor extension of (6) is

$$\mathbf{y}_i \stackrel{\text{ind.}}{\sim} N(\mathbf{f}(x_i), \sigma_\varepsilon^2), \quad (11)$$

where the predictors x_i , $1 \leq i \leq n$, are each 2×1 vectors. In many applications, the x_i s correspond to geographical location but they could be measurements on any pair of predictors for which a bivariate mean function might be entertained. Mixed model-based penalized splines can handle this bivariate predictor case by extending (5) to

$$\mathbf{f}(x) = \beta_0 + \beta_1^\top x + \sum_{k=1}^K u_k z_k(x), \quad u_k \stackrel{\text{ind.}}{\sim} N(0, \tau_u^{-1}) \quad (12)$$

and setting z_k to be appropriate bivariate spline functions. There are several options for doing this (see, e.g., Ruppert *et al.* 2009, Section 2.2). A relatively simple choice is described here and used in the examples. It is based on thin plate spline theory, and corresponds to Section 13.5 of Ruppert *et al.* (2003). The first step is to choose the number K of bivariate knots and their locations. We denote these 2×1 vectors by $\kappa_1, \dots, \kappa_K$. Our default rule for choosing knot locations involves feeding the x_i s and K into the clustering algorithm known as CLARA (Kaufman and Rousseeuw 1990) and setting the κ_k to be cluster centers. Next, form the matrices

$$X = \begin{bmatrix} 1 & x_1^\top \\ \vdots & \vdots \\ 1 & x_n^\top \end{bmatrix},$$

$$Z_K = \begin{bmatrix} \|x_1 - \kappa_1\|^2 \log \|x_1 - \kappa_1\| & \cdots & \|x_1 - \kappa_K\|^2 \log \|x_1 - \kappa_K\| \\ \vdots & \ddots & \vdots \\ \|x_n - \kappa_1\|^2 \log \|x_n - \kappa_1\| & \cdots & \|x_n - \kappa_K\|^2 \log \|x_n - \kappa_K\| \end{bmatrix}, \text{ and}$$

$$\Omega = \begin{bmatrix} \|\kappa_1 - \kappa_1\|^2 \log \|\kappa_1 - \kappa_1\| & \cdots & \|\kappa_1 - \kappa_K\|^2 \log \|\kappa_1 - \kappa_K\| \\ \vdots & \ddots & \vdots \\ \|\kappa_K - \kappa_1\|^2 \log \|\kappa_K - \kappa_1\| & \cdots & \|\kappa_K - \kappa_K\|^2 \log \|\kappa_K - \kappa_K\| \end{bmatrix}.$$

Based on the singular value decomposition $\mathbf{\Omega} = \mathbf{U}\text{diag}(\mathbf{d})\mathbf{V}^\top$, compute $\mathbf{\Omega}^{1/2} = \mathbf{U}\text{diag}(\sqrt{\mathbf{d}})\mathbf{V}^\top$ and set $\mathbf{Z} = \mathbf{Z}_K\mathbf{\Omega}^{-1/2}$. Then

$$z_k(\mathbf{x}_i) = \text{the } (i, k)\text{th entry of } \mathbf{Z} \quad \text{with} \quad \mathbf{u}_k \stackrel{\text{ind.}}{\sim} N(\mathbf{0}, \tau_u^{-1}).$$

2.6. Support files

The code for fitting the semiparametric models considered in this paper are a series of a R and C# source code files utilizing *Infer.NET*. More specifically, for a particular example, an R script exists which reads in or simulates data for the example, constructs relevant basis functions, saves relevant information to comma separated variable (.csv) files, compiles corresponding C# source code files and runs them. The semiparametric model is then specified and fit utilizing *Infer.NET* libraries with the results saved to .csv files. Finally, the *Infer.NET* results are loaded into R and processed. Note that the C# source code files are compiled using Visual Studio 2010 from the command line in the Microsoft Windows operating system.

We have put together a set of files aimed at supporting and easing implementation of the above pipeline. These files include functions for the construction of spline basis functions, and the assessment of the quality of approximated posterior distributions. These support files are summarized below.

InferNETSupport.R This R script loads all of the required packages and support functions used by the examples presented in this paper. The main function in the script is called `RunInferNET()`. This function takes an .cs file (containing *Infer.NET* code), copies this file to a specified support directory renaming it `RunMe.cs`. The file `RunMe.cs` is part of a template C# project called `InferNETSupport.csproj` which is then compiled using the MSBuild program (a part of Microsoft's .NET framework).

DataTable.cs This file defines the 'DataTable' class for reading a data matrix stored in a comma separated file format for use by *Infer.NET*.

SaveData.cs This file defines the 'SaveData' class for taking a data matrix in *Infer.NET* and saving it as a comma separated file.

splineFunctions.R This R script contains a number of functions for constructing various spline basis functions. These include `ZOSull()`, a function which constructs O'Sullivan splines (Wand and Ormerod 2008), `Ztps()`, which constructs thin plate splines (Ruppert et al. 2003), and the function `save.default.knots.2D()`, which selects knots for a two-dimensional dataset (slightly modified from the **SemiPar** package in R).

approximateInferenceFunctions.R This R script contains `accVarApp()` and `summMCMC()`, functions for summarizing the results of variational approximation and MCMC methods. Densities implemented in `accVarApp()` include normal, inverse-gamma, sqrt-inverse-gamma, beta, Bernoulli, discrete, normal mixtures and some non-standard densities considered in Wand, Ormerod, Padoan, and Frühwirth (2011).

mixFunctions.R This R script contains the functions `dnormMix()`, `pnormMix()` and `qnormMix()` for calculating the density, distribution and quantile functions for a univariate normal

mixture. It also contains the function `logOfSum()` for calculating the log of a sum of the exponential of a vector in a numerically safe manner.

All of the semiparametric regression models considered in this paper are implemented using our pipeline in the `InferNetScripts` directly from the `.zip` file containing the supplementary material to this paper. We anticipate that implementing code for semiparametric regression models similar to those considered in this paper should be relatively straightforward by modifying the code in these files directly.

3. Simple semiparametric model

The first example involves the simple semiparametric regression model

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \sum_{k=1}^K u_k z_k(x_{2i}) + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{ind.}}{\sim} N(0, \tau_\varepsilon^{-1}), \quad 1 \leq i \leq n, \quad (13)$$

where $z_k(\cdot)$ is a set of spline basis functions as described in Section 2.5. The corresponding Bayesian mixed model can then be represented by

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}\mathbf{I}), \quad \mathbf{u} | \tau_u \sim N(\mathbf{0}, \tau_u^{-1}\mathbf{I}), \quad (14)$$

$$\boldsymbol{\beta} \sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), \quad \tau_u^{-1/2} \sim \text{Half-Cauchy}(\mathbf{A}_u), \quad \tau_\varepsilon^{-1/2} \sim \text{Half-Cauchy}(\mathbf{A}_\varepsilon),$$

where τ_β , \mathbf{A}_u and \mathbf{A}_ε are user-specified hyperparameters and

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_K \end{bmatrix},$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{1} & \mathbf{x}_{11} & \mathbf{x}_{21} \\ \vdots & \vdots & \\ \mathbf{1} & \mathbf{x}_{1n} & \mathbf{x}_{2n} \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} z_1(\mathbf{x}_{21}) & \cdots & z_K(\mathbf{x}_{21}) \\ \vdots & \ddots & \vdots \\ z_1(\mathbf{x}_{2n}) & \cdots & z_K(\mathbf{x}_{2n}) \end{bmatrix}.$$

Infer.NET 2.6 does not support direct fitting of model (14) under the product restriction

$$q(\boldsymbol{\beta}, \mathbf{u}, \tau_u, \tau_\varepsilon) = q(\boldsymbol{\beta}, \mathbf{u}) q(\tau_u, \tau_\varepsilon). \quad (15)$$

We get around this by employing the same trick as described in Section 3.2 of Wang and Wand (2011). It entails the introduction of the auxiliary $n \times 1$ data vector \mathbf{a} . By setting the observed data for \mathbf{a} equal to 0 and by assuming a very small number for κ , fitting the model in (16) provides essentially the same result as fitting the model in (14). The actual model implemented in Infer.NET is

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}\mathbf{I}), \quad \mathbf{a} | \boldsymbol{\beta}, \mathbf{u}, \tau_u \sim N\left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \begin{bmatrix} \tau_\beta^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tau_u^{-1}\mathbf{I} \end{bmatrix}\right),$$

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \sim N(\mathbf{0}, \kappa^{-1}\mathbf{I}), \quad \tau_u | \mathbf{b}_u \sim \text{Gamma}(\tfrac{1}{2}, \mathbf{b}_u),$$

$$\mathbf{b}_u \sim \text{Gamma}(\tfrac{1}{2}, 1/A_u^2), \quad \tau_\varepsilon | \mathbf{b}_\varepsilon \sim \text{Gamma}(\tfrac{1}{2}, \mathbf{b}_\varepsilon), \quad \mathbf{b}_\varepsilon \sim \text{Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2), \quad (16)$$

with the inputted a vector containing zeros. While the full *Infer.NET* code is included in Appendix A, we will confine discussion here to the key parts of the code that specify (16). Specification of the prior distributions for τ_u and τ_ε can be achieved using the following *Infer.NET* code:

```
Variable<double> tauu = Variable.GammaFromShapeAndRate(0.5,
  Variable.GammaFromShapeAndRate(0.5, Math.Pow(Au, -2))).Named("tauu");
Variable<double> tauEps = Variable.GammaFromShapeAndRate(0.5,
  Variable.GammaFromShapeAndRate(0.5, Math.Pow(Aeps, -2))).Named("tauEps");
```

(17)

while the likelihood in (16) is specified by the following line:

```
y[index] = Variable.GaussianFromMeanAndPrecision(
  Variable.InnerProduct(betauWork, cvec[index]), tauEps);
```

(18)

The variable `index` represents a range of integers from 1 to n and enables loop-type structures. Finally, the inference engine is specified to be variational message passing via the code:

```
InferenceEngine engine = new InferenceEngine();
engine.Algorithm = new VariationalMessagePassing();
engine.NumberOfIterations = nIterVB;
```

(19)

with `nIterVB` denoting the number of mean field variational Bayes iterations. Note that this *Infer.NET* code is treating the coefficient vector

$$\begin{bmatrix} \beta \\ u \end{bmatrix}$$

as an entity, in keeping with product restriction (15).

Figures 2 and 3 summarize the results from *Infer.NET* fitting of (16) to a data set on the yields (g/plant) of 84 white Spanish onions crops in two locations: Purnong Landing and Virginia, South Australia. The response variable, y is the logarithm of yield, whilst the predictors are indicator of location being Virginia (x_1) and areal density of the plants (plants/m²) (x_2). The hyperparameters were set at $\tau_\beta = 10^{-10}$, $A_\varepsilon = 10^5$, $A_u = 10^5$ and $\kappa = 10^{-10}$, while the number of mean field variational Bayes iterations was fixed at 100.

As a benchmark, Bayesian inference via MCMC was performed. To this end, the following BUGS program was used:

```
for (i in 1:n) {
  mu[i] <- (beta0 + beta1 * x1[i] + beta2 * x2[i] + inprod(u[], Z[i, ]))
  y[i] ~ dnorm(mu[i], tauEps)
}
for (k in 1:K) {
  u[k] ~ dnorm(0, tauu)
}
beta0 ~ dnorm(0, 1.0E-10); beta1 ~ dnorm(0, 1.0E-10)
beta2 ~ dnorm(0, 1.0E-10);
bu ~ dgamma(0.5, 1.0E-10); tauu ~ dgamma(0.5, bu)
bEps ~ dgamma(0.5, 1.0E-10); tauEps ~ dgamma(0.5, bEps)
```

(20)

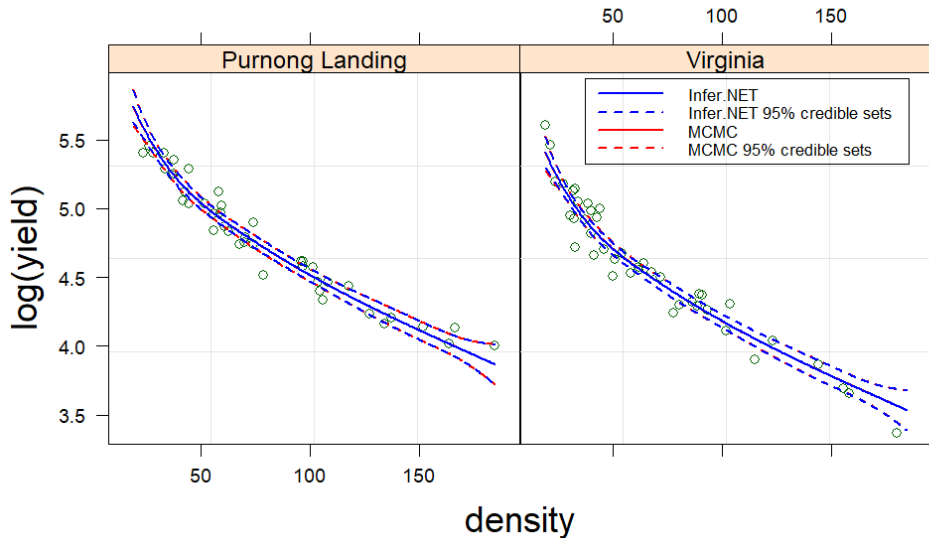


Figure 2: Onions data. Fitted regression line and pointwise 95% credible intervals for variational Bayesian inference by Infer.NET and MCMC.

Parameter (θ)	MCMC		Infer.NET	
	$E(\theta y)$	95% credible interval	$E(\theta y)$	95% credible interval
β_0	0.512	(0.416, 0.606)	0.516	(0.428, 0.604)
β_1	-0.716	(-0.817, -0.614)	-0.717	(-0.817, -0.618)
β_2	-0.942	(-1.003, -0.882)	-0.944	(-1.001, -0.886)
τ_ε	18.9	(13.3, 25.4)	19.0	(13.7, 25.1)
τ_u	0.598	(0.008, 3.332)	0.316	(0.038, 0.881)

Table 2: Onions data. Parameter point estimates and 95% credible intervals for variational Bayesian inference by Infer.NET and MCMC.

A burn-in length of 50000 was used, while 1 million samples were obtained from the posterior distributions. Finally, a thinning factor of 5 was used. The posterior densities for MCMC were produced based on kernel density estimation with plug-in bandwidth selection via the R package **KernSmooth** (Wand 2015). Figure 2 displays the fitted regression lines and pointwise 95% credible intervals. The estimated regression lines and credible intervals from Infer.NET and MCMC fitting are highly similar. Figure 3 visualizes the approximate posterior density functions for β_1 , τ_ε^{-1} and the effective degrees of freedom of the fit. The variational Bayes approximations for β_1 and τ_ε^{-1} are rather accurate. The approximate posterior density function for the effective degrees of freedom for variational Bayesian inference was obtained based on Monte Carlo samples of size 1 million from the approximate posterior distributions of τ_ε^{-1} and τ_u^{-1} .

A summary of the parameter point estimates and 95% credible intervals for this model is given in Table 2. In this table we see that inferences for Infer.NET and MCMC are identical to the first two significant figures for β_0 , β_1 , β_2 and τ_ε . However, τ_u is underestimated by Infer.NET with credible intervals which are too small.

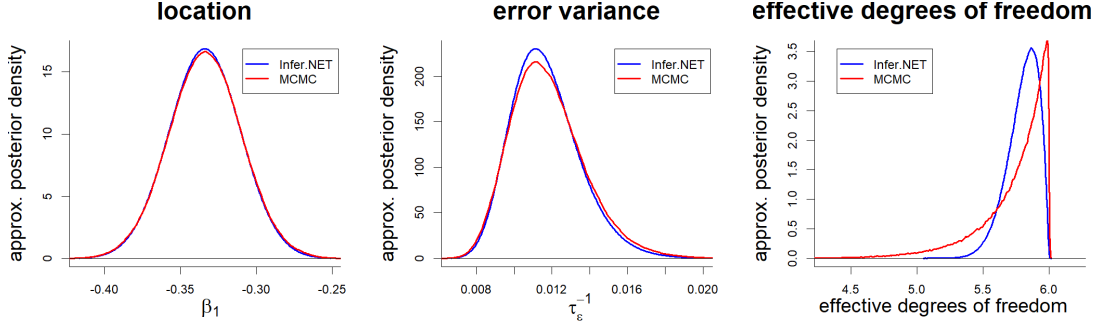


Figure 3: Variational Bayes approximate posterior density functions produced by Infer.NET and MCMC posterior density functions of key model parameters for fitting a simple semi-parametric model to the onions data set.

4. Generalized additive model

The next example illustrates binary response variable regression through the model

$$\begin{aligned} \mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli}(F(\{\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}\}_i)), & \mathbf{u} | \tau_u &\sim N(\mathbf{0}, \tau_u^{-1}\mathbf{I}), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), & \tau_u^{-1/2} &\sim \text{Half-Cauchy}(\mathbf{A}_u), \quad 1 \leq i \leq n, \end{aligned} \quad (21)$$

where $F(\cdot)$ denotes an inverse link function and \mathbf{X} , \mathbf{Z} , τ_β and \mathbf{A}_u are defined as in the previous section. Typical choices of $F(\cdot)$ correspond to logistic regression and probit regression. As before, introduction of the auxiliary variables $a \equiv 0$ enables Infer.NET fitting of the Bayesian mixed model

$$\begin{aligned} \mathbf{y}_i | \boldsymbol{\beta}, \mathbf{u} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli}(F(\{\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}\}_i)), \\ \mathbf{a} | \boldsymbol{\beta}, \mathbf{u}, \tau_u &\sim N\left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \begin{bmatrix} \tau_\beta^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tau_u^{-1}\mathbf{I} \end{bmatrix}\right), \\ \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} &\sim N(\mathbf{0}, \boldsymbol{\kappa}^{-1}\mathbf{I}), & \tau_u | \mathbf{b}_u &\sim \text{Gamma}(\tfrac{1}{2}, \mathbf{b}_u), \\ \mathbf{b}_u &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_u^2), & 1 \leq i &\leq n. \end{aligned} \quad (22)$$

The likelihood for the logistic regression case is specified as follows

```
VariableArray<bool> y = Variable.Array<bool>(index).Named("y");
y[index] = Variable.BernoulliFromLogOdds(
  Variable.InnerProduct(betauWork, cvec[index]));
```

(23)

while variational message passing is used for fitting purposes. Setting up the prior for τ_u and specifying the inference engine is done as in code chunk (17) and (19), respectively. For probit regression, the last line of (23) is replaced by

```
y[index] = Variable.IsPositive(Variable.GaussianFromMeanAndVariance(
  Variable.InnerProduct(betauWork, cvec[index]), 1));
```

(24)

and expectation propagation is used

```
engine.Algorithm = new ExpectationPropagation();
```

 (25)

Instead of the half-Cauchy prior in (21), a gamma prior with shape and rate equal to 2 is used for τ_u in the probit regression model

```
Variable<double> tauU = Variable.GammaFromShapeAndRate(2, 2);
```

 (26)

Figure 4 visualizes the results for Infer.NET and MCMC fitting of a straightforward extension of model (22) with inverse-logit and probit link functions to a breast cancer data set (Haberman 1976). Here our MCMC scheme used a burn-in length of 5000 and obtained a further 5000 samples to be used for inference. This data set contains 306 cases from a study that was conducted between 1958 and 1970 at the University of Chicago’s Billings Hospital on the survival of patients who had undergone surgery for breast cancer. The binary response variable represents whether the patient died within 5 years after operation (`died`), while 3 predictor variables are used: age of patient at the time of operation (`age`), year of operation (`year`) and number of positive axillary nodes (`nodes`) detected. The actual model is

$$\text{died}_i | \beta, \mathbf{u} \stackrel{\text{ind.}}{\sim} \text{Bernoulli}(\mathbf{F}(\beta_0 + \mathbf{f}_1(\text{age}_i) + \mathbf{f}_2(\text{year}_i) + \mathbf{f}_3(\text{nodes}_i))), \quad 1 \leq i \leq 306.$$

Note that all predictor variables were standardized prior to model fitting and the following hyperparameters were used: $\tau_\beta = 10^{-10}$, $A_u = 10^5$, $\kappa = 10^{-10}$ and the number of variational Bayes iterations was 100. Figure 4 illustrates that there is good agreement between the results from Infer.NET and MCMC.

5. Robust nonparametric regression with the t distribution

A popular model-based approach for robust regression is to model the response variable to have a t distribution. Outliers occur with moderate probability for low values of the t distribution’s degrees of freedom parameter (Lange, Little, and Taylor 1989). More recently, Staudenmayer, Lake, and Wand (2009) proposed a penalized spline mixed model approach to nonparametric regression using the t distribution.

The robust nonparametric regression model that we consider here is a Bayesian variant of that treated by Staudenmayer *et al.* (2009):

$$\begin{aligned} y_i | \beta_0, \beta_1, \mathbf{u}, \tau_\varepsilon, \nu &\stackrel{\text{ind.}}{\sim} t\left(\beta_0 + \beta_1 x_i + \sum_{k=1}^K u_k z_k(x_i), \tau_\varepsilon^{-1}, \nu\right), \quad 1 \leq i \leq n, \\ \mathbf{u} | \tau_u &\sim \mathbf{N}(\mathbf{0}, \tau_u^{-1} \mathbf{I}), \quad \beta \sim \mathbf{N}(\mathbf{0}, \tau_\beta^{-1} \mathbf{I}), \\ \tau_u^{-1/2} &\sim \text{Half-Cauchy}(A_u), \quad \tau_\varepsilon^{-1/2} \sim \text{Half-Cauchy}(A_\varepsilon), \\ p(\nu) &\text{ discrete on a finite set } N. \end{aligned}$$
 (27)

Restricting the prior distribution of ν to be that of a discrete random variable allows Infer.NET fitting of the model using structured mean field variational Bayes (Saul and Jordan 1996). This extension of ordinary mean field variational Bayes is described in Section 3.1 of Wand *et al.* (2011). Since variational message passing is a special case of mean field variational Bayes this extension also applies. Model (27) can be fitted through calls to Infer.NET with ν fixed at each value in N . The results of each of these fits are combined afterwards. Details are given below.

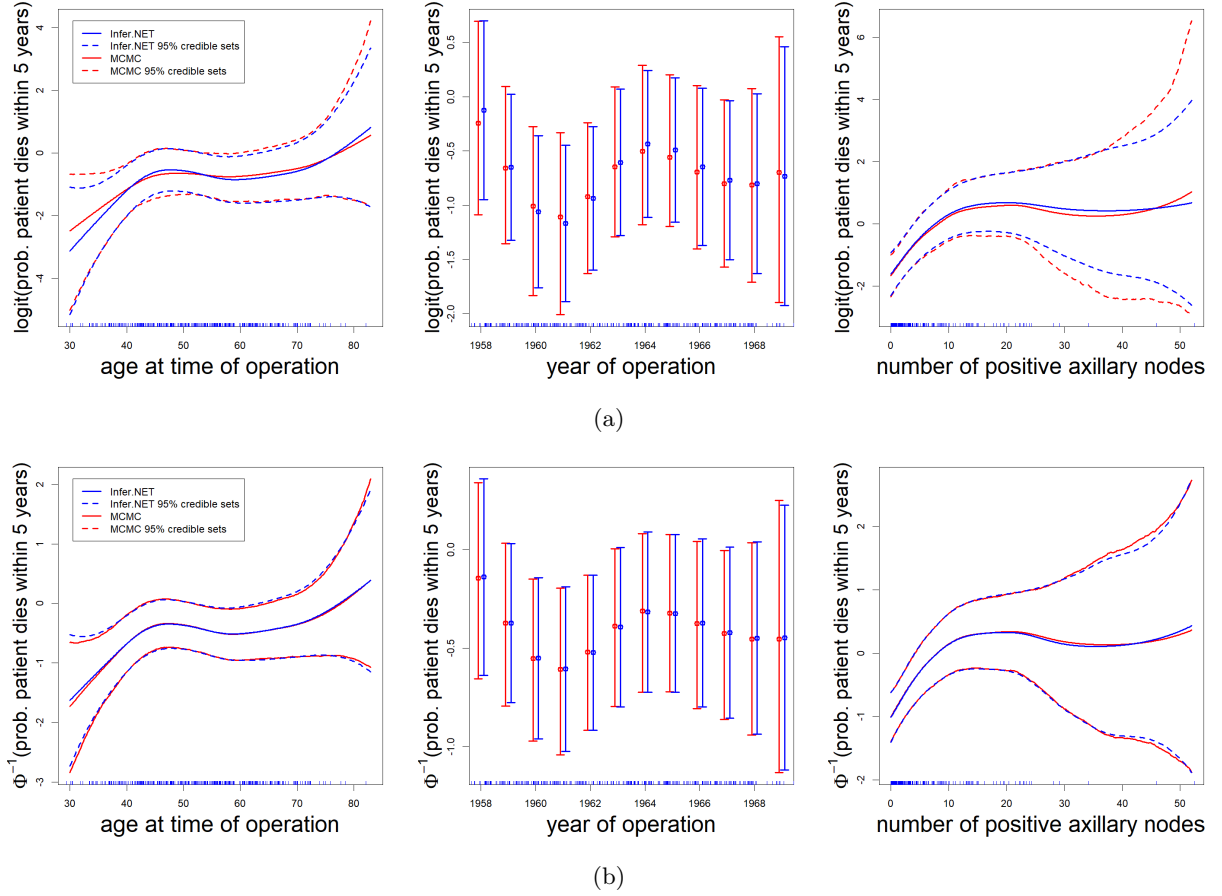


Figure 4: Breast cancer data. Fitted regression lines and pointwise 95% credible intervals for logit (a) and probit (b) regression via variational Bayesian inference by *Infer.NET* and MCMC.

Another challenge concerning (27) is that *Infer.NET* does not support t distribution specifications. We get around this by resorting to the result:

$$\text{If } x|g \sim N(\mu, (g\tau)^{-1}) \quad \text{and} \quad g \sim \text{Gamma}(\frac{\nu}{2}, \frac{\nu}{2}), \quad \text{then } x \sim t(\mu, \tau^{-1}, \nu). \quad (28)$$

As before, we use the $a \equiv 0$ auxiliary data trick described in Section 3 and the half-Cauchy representation (1). These lead to the following suite of models, for each fixed $\nu \in N$, needed to be run in *Infer.NET*:

$$\begin{aligned} y | \beta, \mathbf{u}, \tau_\varepsilon &\sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1} \text{diag}(\mathbf{1}/g)), \\ a | \beta, \mathbf{u}, \tau_u &\sim N\left(\begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix}, \begin{bmatrix} \tau_\beta^{-1} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tau_u^{-1} \mathbf{I} \end{bmatrix}\right), \quad \begin{bmatrix} \beta \\ \mathbf{u} \end{bmatrix} \sim N(\mathbf{0}, \kappa^{-1} \mathbf{I}) \\ \mathbf{u} | \tau_u &\sim N(\mathbf{0}, \tau_u^{-1} \mathbf{I}), \quad g_i | \nu \stackrel{\text{ind.}}{\sim} \text{Gamma}(\frac{\nu}{2}, \frac{\nu}{2}), \quad \beta \sim N(\mathbf{0}, \tau_\beta^{-1} \mathbf{I}), \quad (29) \\ \tau_u | b_u &\sim \text{Gamma}(\frac{1}{2}, b_u), \quad b_u \sim \text{Gamma}(\frac{1}{2}, 1/A_u^2), \\ \tau_\varepsilon | b_\varepsilon &\sim \text{Gamma}(\frac{1}{2}, b_\varepsilon), \quad b_\varepsilon \sim \text{Gamma}(\frac{1}{2}, 1/A_\varepsilon^2), \end{aligned}$$

where the matrix notation of (8) and (9) is used, $g = [g_1, \dots, g_n]^\top$ and τ_β , A_u and A_ε are user-specified hyperparameters with default values $\tau_\beta = 10^{-10}$, $A_u = A_\varepsilon = 10^5$. The prior for ν was set to be a uniform distribution over the atom set N , which is set equal to 30 equally-spaced numbers between 0.05 and 10 inclusive.

After obtaining fits of (29) for each $\nu \in N$, the approximate posterior densities are obtained from

$$\begin{aligned} \mathbf{q}(\beta, \mathbf{u}) &= \sum_{\nu \in N} \mathbf{q}_\nu(\nu) \mathbf{q}(\beta, \mathbf{u} | \nu), & \mathbf{q}(\tau_u) &= \sum_{\nu \in N} \mathbf{q}_\nu(\nu) \mathbf{q}(\tau_u | \nu), \\ \mathbf{q}(\tau_\varepsilon) &= \sum_{\nu \in N} \mathbf{q}(\nu) \mathbf{q}(\tau_\varepsilon | \nu), & \text{with } \mathbf{q}(\nu) &= \underline{p}(\nu | \mathbf{y}) = \frac{\mathbf{p}(\nu) \underline{p}(\mathbf{y} | \nu)}{\sum_{\nu' \in N} \mathbf{p}(\nu') \underline{p}(\mathbf{y} | \nu')}, \end{aligned}$$

where $\underline{p}(y | \nu)$ denotes the variational lower bound on the conditional likelihood $p(y | \nu)$.

Specification of the prior distribution for g is achieved using the following Infer.NET code:

```
VariableArray<double> g = Variable.Array<double>(index);
g[index] = Variable.GammaFromShapeAndRate(nu/2, 2/nu).ForEach(index);
```

 (30)

while the likelihood in (29) is specified by:

```
y[index] = Variable.GaussianFromMeanAndPrecision(
  Variable.InnerProduct(betauWork, cvec[index]), g[index] * tauEps);
```

 (31)

The lower bound $\log \underline{p}(y | \nu)$ can be obtained by creating a mixture of the current model with an empty model in Infer.NET. The learned mixing weight is then equal to the marginal log-likelihood. Therefore, an auxiliary Bernoulli variable is set up:

```
Variable<bool> auxML = Variable.Bernoulli(0.5).Named("auxML");
IfBlock model = Variable.If(auxML);
```

 (32)

The normal code for fitting the model in (29) is then enclosed with

```
IfBlock model = Variable.If(auxML);
```

 (33)

and

```
model.CloseBlock();
```

 (34)

Finally, the lower bound $\log \underline{p}(y | \nu)$ is obtained from:

```
double marginalLogLikelihood = engine.Infer<Bernoulli>(auxML).LogOdds;
```

 (35)

Figure 5 presents the results of the structured mean field variational Bayes analysis using Infer.NET fitting of model (29) to a data set on a respiratory experiment conducted by Professor Russ Hauser at Harvard School of Public Health, Boston, USA. The data correspond to 60 measurements on one subject during two separate respiratory experiments. The response variable y_i represents the log of the adjusted time of exhalation for x_i equal to the time in seconds since exposure to air containing a particulate matter. The adjusted time of exhalation is obtained by subtracting the average time of exhalation at baseline, prior to exposure to filtered air. Interest centers upon the mean response as a function of time. The predictor and response variable were both standardized prior to Infer.NET analysis. The following hyperparameters were chosen: $\tau_\beta = 10^{-10}$, $A_\varepsilon = 10^5$, $A_u = 10^5$ and $\kappa = 10^{-10}$, while the number

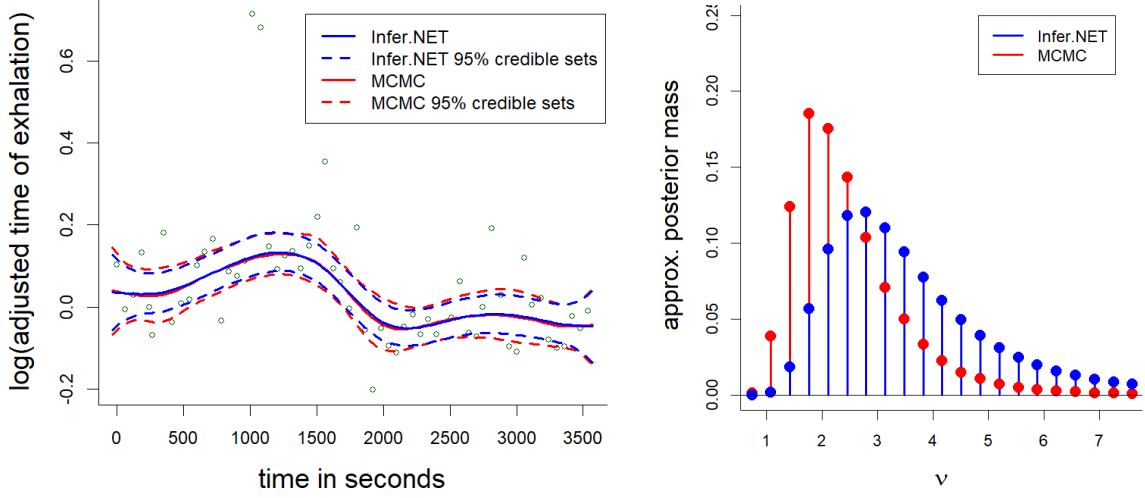


Figure 5: Structured mean field variational Bayesian, based on *Infer.NET*, and MCMC fitting of the robust nonparametric regression model (27) to the respiratory experiment data. Left: Posterior mean and pointwise 95% credible sets for the regression function. Right: Approximate posterior function for the degrees of freedom parameter ν .

of variational Bayes iterations equaled 100. Bayesian inference via MCMC was performed as a benchmark where a burn-in length of 5000 was used and a further 5000 samples were obtained for inference.

Figure 5 shows that the variational Bayes fit and pointwise 95% credible sets are close to the ones obtained using MCMC. Finally, the approximate posterior probability function and the posterior from MCMC for the degrees of freedom ν are compared. The *Infer.NET* and MCMC results coincide quite closely.

6. Semiparametric mixed model

Since semiparametric regression models based on penalized splines fit in the mixed model framework, semiparametric longitudinal data analysis can be performed by fusion with classical mixed models (Ruppert *et al.* 2003). In this section we illustrate the use of *Infer.NET* for fitting the class of semiparametric mixed models having the form:

$$y_{ij} | \beta, \mathbf{u}, U_i, \tau_\varepsilon \stackrel{\text{ind.}}{\sim} N \left(\beta_0 + \beta^\top x_{ij} + U_i + \sum_{k=1}^K u_k z_k(s_{ij}), \tau_\varepsilon^{-1} \right), \quad (36)$$

$$U_i | \tau_U \stackrel{\text{ind.}}{\sim} N(0, \tau_U^{-1}), \quad 1 \leq j \leq n_i, \quad 1 \leq i \leq m,$$

for analysis of longitudinal data sets such as the one described in Bachrach, Hastie, Wang, Narasimhan, and Marcus (1999). Here x_{ij} is a vector of predictors that enter the model linearly and s_{ij} is another predictor that enters the model nonlinearly via penalized splines. For each $1 \leq i \leq m$, U_i denotes the random intercept for the i th subject. The corresponding

Bayesian mixed model is represented by

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}\mathbf{I}), & \mathbf{u} | \tau_U, \tau_u &\sim N\left(\mathbf{0}, \begin{bmatrix} \tau_U^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \tau_u^{-1}\mathbf{I} \end{bmatrix}\right), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), & \tau_u^{-1/2} &\sim \text{Half-Cauchy}(A_u), & \tau_\varepsilon^{-1/2} &\sim \text{Half-Cauchy}(A_\varepsilon), \\ & & \tau_U^{-1/2} &\sim \text{Half-Cauchy}(A_U), \end{aligned} \quad (37)$$

where \mathbf{y} , $\boldsymbol{\beta}$ and \mathbf{X} are defined in a similar manner as in the previous sections and τ_β , A_u , A_ε and A_U are user-specified hyperparameters. Introduction of the random intercepts results in

$$\mathbf{u} = \begin{bmatrix} \mathbf{U}_1 \\ \vdots \\ \mathbf{U}_m \\ \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_K \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{1} & \cdots & \mathbf{0} & z_1(\mathbf{s}_{11}) & \cdots & z_K(\mathbf{s}_{11}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1} & \cdots & \mathbf{0} & z_1(\mathbf{s}_{1n_1}) & \cdots & z_K(\mathbf{s}_{1n_1}) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1} & z_1(\mathbf{s}_{m1}) & \cdots & z_K(\mathbf{s}_{m1}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1} & z_1(\mathbf{s}_{mn_m}) & \cdots & z_K(\mathbf{s}_{mn_m}) \end{bmatrix}.$$

We again use the auxiliary data vector $\mathbf{a} \equiv \mathbf{0}$ to allow direct Infer.NET fitting of the following Bayesian longitudinal penalized spline model:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}\mathbf{I}), & \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} &\sim N(\mathbf{0}, \kappa^{-1}\mathbf{I}), \\ \mathbf{a} | \boldsymbol{\beta}, \mathbf{u}, \tau_U, \tau_u &\sim N\left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \begin{bmatrix} \tau_\beta^{-1}\mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tau_U^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \tau_u^{-1}\mathbf{I} \end{bmatrix}\right), \end{aligned} \quad (38)$$

$$\tau_u | b_u \sim \text{Gamma}(\tfrac{1}{2}, b_u), \quad b_u \sim \text{Gamma}(\tfrac{1}{2}, 1/A_u^2), \quad \tau_\varepsilon | b_\varepsilon \sim \text{Gamma}(\tfrac{1}{2}, 1/b_\varepsilon),$$

$$b_\varepsilon \sim \text{Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2), \quad \tau_U | b_U \sim \text{Gamma}(\tfrac{1}{2}, b_U), \quad b_U \sim \text{Gamma}(\tfrac{1}{2}, 1/A_U^2).$$

Specification of the prior distribution for τ_U can be achieved in Infer.NET in a similar manner as the prior specification in code chunk (17).

Figure 6 shows the Infer.NET fits of (38) to the spinal bone mineral density data (Bachrach *et al.* 1999). A population of 230 female subjects aged between 8 and 27 was followed over time and each subject contributed either one, two, three, or four spinal bone mineral density measurements. Age enters the model nonlinearly and corresponds to s_{ij} in (36). Data on ethnicity are available and the entries of x_{ij} correspond to the indicator variables for Black (x_{1ij}), Hispanic (x_{2ij}) and White (x_{3ij}), with Asian ethnicity corresponding to the baseline. The following hyperparameters were chosen: $\tau_\beta = 10^{-10}$, $A_\varepsilon = 10^5$, $A_u = 10^5$, $A_U = 10^5$ and $\kappa = 10^{-10}$, while the number of mean field variational Bayes iterations was set to 100. The MCMC fits used 5000 burn-in samples and 100,000 samples for inference. Figure 6 suggests that there exists a statistically significant difference in mean spinal bone mineral density

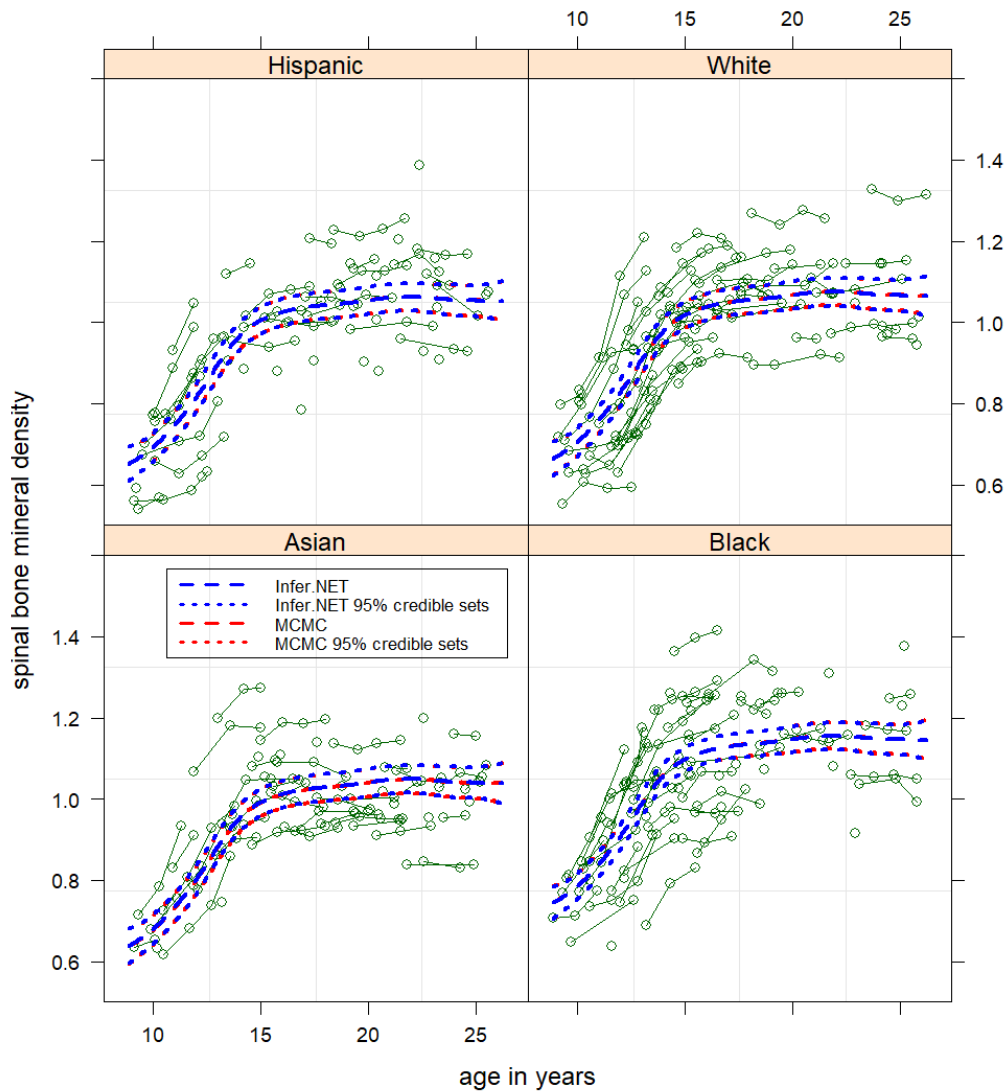


Figure 6: Spinal bone mineral density data. Fitted regression line and pointwise 95% credible intervals based on mean field variational Bayesian inference by *Infer.NET* and MCMC. The regression lines represent the smoothed fits for each group in the longitudinal study. Note that the fitted regression lines for both *Infer.NET* and MCMC are nearly identical.

between Asian and Black subjects. This difference is confirmed by the approximate posterior density functions in Figure 7. No statistically significant difference is found between Asian and Hispanic subjects and between Asian and White subjects.

7. Geoadditive model

We turn our attention to geostatistical data, for which the response variable is geographically referenced and the extension of generalized additive models, known as geoadditive models ([Kammann and Wand 2003](#)), applies. Such models allow for a pair of continuous predictors, typically geographical location, to have a bivariate functional impact on the mean response.

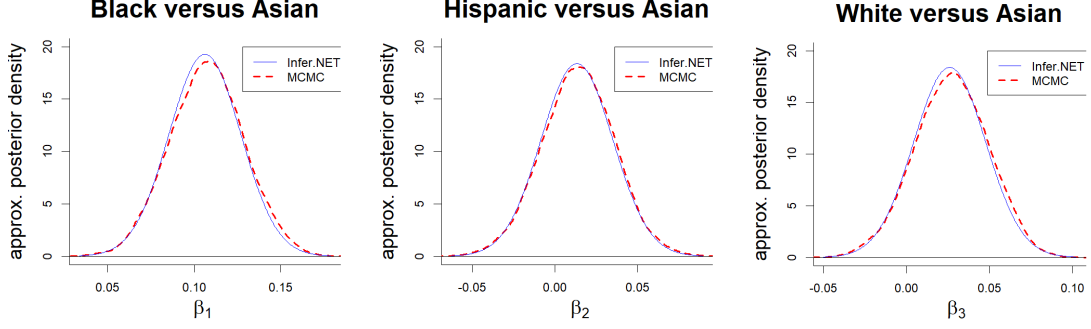


Figure 7: Variational Bayes approximate posterior density functions produced by Infer.NET and MCMC of the ethnic group parameters when fitting a simple semiparametric mixed model to the spinal bone mineral density data set.

An example geoaddivitive model is given by

$$y_i \sim N\left(f_1(x_{1i}) + f_2(x_{2i}) + f_{\text{geo}}(x_{3i}, x_{4i}), \tau_\varepsilon^{-1}\right), \quad 1 \leq i \leq n. \quad (39)$$

It can be handled using the spline basis described in Section 2.5 as follows:

$$y_i | \beta, u_1, u_2, u^{\text{geo}}, \tau_\varepsilon \stackrel{\text{ind.}}{\sim} N\left(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 x_{4i} + \sum_{k=1}^{K_1} u_{1k} z_{1k}(x_{1i}) + \sum_{k=1}^{K_2} u_{2k} z_{2k}(x_{2i}) + \sum_{k=1}^{K_{\text{geo}}} u_k^{\text{geo}} z_k^{\text{geo}}(x_{3i}, x_{4i}), \tau_\varepsilon^{-1}\right), \quad 1 \leq i \leq n, \quad (40)$$

where the z_{1k} and z_{2k} are univariate spline basis functions as used in each of this article's previous examples and the z_k^{geo} are the bivariate spline basis functions, described in Section 2.5. The corresponding Bayesian mixed model is represented by

$$y | \beta, u, \tau_\varepsilon \sim N(X\beta + Zu, \tau_\varepsilon^{-1}I), \quad \beta \sim N(0, \tau_\beta^{-1}I),$$

$$u = \begin{bmatrix} u_1 \\ u_2 \\ u^{\text{geo}} \end{bmatrix} \Bigg| \tau_{u1}, \tau_{u2}, \tau_{\text{geo}} \sim N\left(0, \begin{bmatrix} \tau_{u1}^{-1}I & 0 & 0 \\ 0 & \tau_{u2}^{-1}I & 0 \\ 0 & 0 & \tau_{\text{geo}}^{-1}I \end{bmatrix}\right), \quad (41)$$

$$\tau_u^{-1/2} \sim \text{Half-Cauchy}(A_u), \quad \tau_\varepsilon^{-1/2} \sim \text{Half-Cauchy}(A_\varepsilon),$$

$$\tau_{\text{geo}}^{-1/2} \sim \text{Half-Cauchy}(A_{\text{geo}}),$$

where y, β, u and X are defined in a similar manner as in the previous sections and $\tau_\beta, A_u, A_\varepsilon$ and A_{geo} are user-specified hyperparameters. The Z matrix is

$$Z = \begin{bmatrix} z_{11}(x_{11}) & \cdots & z_{1K_1}(x_{11}) & z_{21}(x_{21}) & \cdots & z_{2K_2}(x_{21}) & z_1^{\text{geo}}(x_{31}, x_{41}) & \cdots & z_{K_{\text{geo}}}^{\text{geo}}(x_{31}, x_{41}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ z_{11}(x_{1n}) & \cdots & z_{1K_1}(x_{1n}) & z_{21}(x_{2n}) & \cdots & z_{2K_2}(x_{2n}) & z_1^{\text{geo}}(x_{3n}, x_{4n}) & \cdots & z_{K_{\text{geo}}}^{\text{geo}}(x_{3n}, x_{4n}) \end{bmatrix}.$$

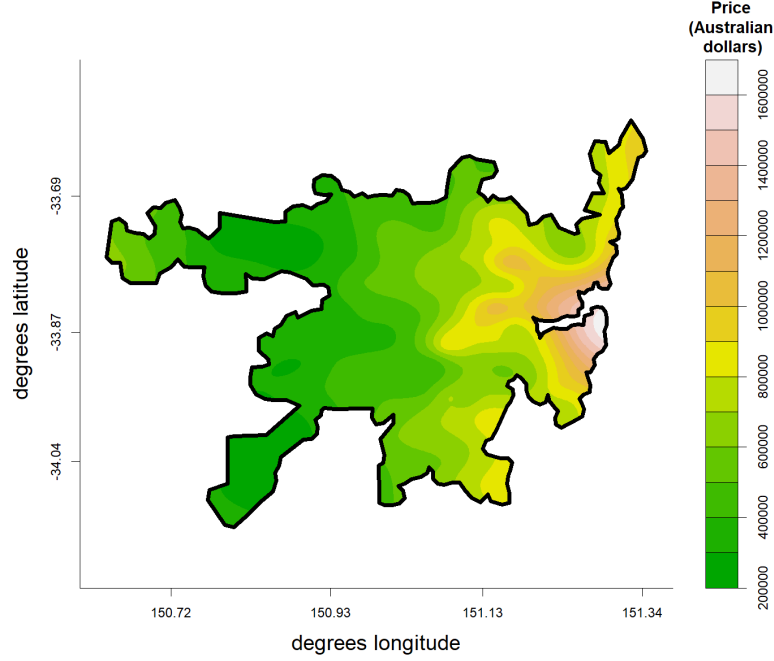


Figure 8: Sydney real estate data. Estimated housing prices for the Sydney metropolitan area for variational Bayesian inference by *Infer.NET*.

Infer.NET can be used to fit the Bayesian geoaddivitive model using

$$\begin{aligned}
 \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \tau_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \tau_\varepsilon^{-1}\mathbf{I}), & \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} &\sim N(\mathbf{0}, \boldsymbol{\kappa}^{-1}\mathbf{I}), \\
 \mathbf{a} | \boldsymbol{\beta}, \mathbf{u}, \tau_u, \tau_{\text{geo}} &\sim N\left(\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix}, \begin{bmatrix} \tau_\beta^{-1}\mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tau_{u1}^{-1}\mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \tau_{u2}^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \tau_{\text{geo}}^{-1}\mathbf{I} \end{bmatrix}\right), \\
 \tau_{u1} | b_{u1} &\sim \text{Gamma}(\tfrac{1}{2}, b_{u1}), & b_{u1} &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_{u1}^2), & \tau_{u2} | b_{u2} &\sim \text{Gamma}(\tfrac{1}{2}, b_{u2}), \\
 b_{u2} &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_{u2}^2), & \tau_\varepsilon | b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, b_\varepsilon), & b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2), \\
 \tau_{\text{geo}} | b_{\text{geo}} &\sim \text{Gamma}(\tfrac{1}{2}, b_{\text{geo}}), & b_{\text{geo}} &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_{\text{geo}}^2)
 \end{aligned}$$

where $\mathbf{a} \equiv \mathbf{0}$ as in all previous examples.

We illustrate geoaddivitive model fitting in *Infer.NET* using data on residential property prices of 37,676 residential properties that were sold in Sydney, Australia, during 2001. The data were assembled as part of an unpublished study by A. Chernih and M. Sherris at the University of New South Wales, Australia. The response variable is the logarithm of the sale price in Australian dollars. Apart from geographical location, several predictor variables are available. For this example we selected weekly income in Australian dollars, the distance from the coastline in kilometers, the particulate matter 10 level and the nitrogen dioxide level. The model is a straightforward extension of the geoaddivitive model conveyed by (39)–(41). In addition, all predictor variables and the dependent variable were standardized before model fitting. The hyperparameters are set to $\tau_\beta = 10^{-10}$, $A_\varepsilon = 10^5 = A_{u_1} = \dots = A_{u_4} = A_{\text{geo}} =$

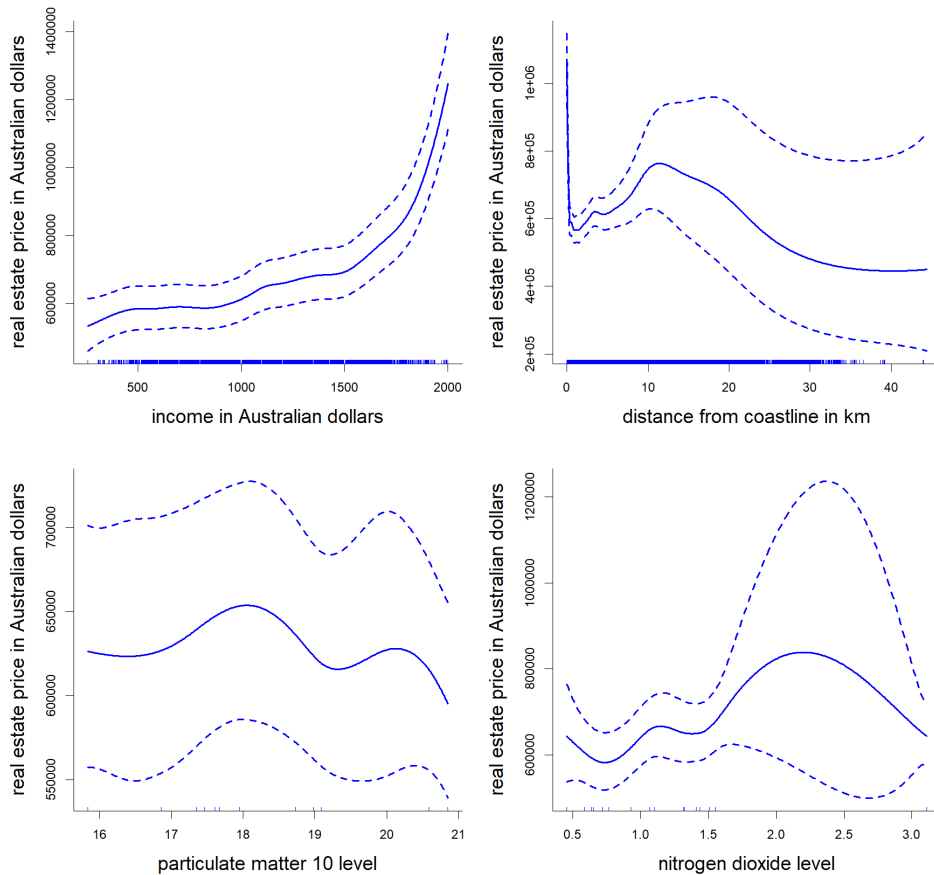


Figure 9: Sydney real estate data. Fitted regression line and pointwise 95% credible intervals for variational Bayesian inference by Infer.NET.

10^5 , while the number of variational message passing iterations was equal to 100 and κ set to 10^{-10} .

Figure 8 summarizes the housing prices for the Sydney metropolitan area based on the fitted Infer.NET model, while fixing the four predictor variables at their mean levels. This result clearly shows that the sale price is higher for houses in Sydney’s Eastern and Northern suburbs whereas it is strongly decreased for houses in Sydney’s Western suburbs.

Figure 9 visualizes the fitted regression line and pointwise 95% credible intervals for income, distance from the coastline, particulate matter 10 level and nitrogen dioxide level at a fixed geographical location (longitude = 151.10° , latitude = -33.91°). As expected, a higher income is associated with a higher sale price and houses close to the coastline are more expensive. The sale price is lower for a higher particulate matter 10 level, while a lower nitrogen dioxide level is generally associated with a lower sale price.

8. Bayesian lasso regression

A Bayesian approach to lasso (least absolute shrinkage selection operator) regression was proposed by Park and Casella (2008). For linear regression, the lasso approach essentially

involves replacing

$$\beta_j | \tau_\beta \stackrel{\text{ind.}}{\sim} N(0, \tau_\beta^{-1}) \quad \text{by} \quad \beta_j | \tau_\beta \stackrel{\text{ind.}}{\sim} \text{Laplace}(0, \tau_\beta^{-1}), \quad (42)$$

where β_j is the coefficient of the j th predictor variable. Replacement (42) corresponds to changing the form of the penalty from

$$\lambda \sum_{j=1}^p \beta_j^2 \quad \text{to} \quad \lambda \sum_{j=1}^p |\beta_j|.$$

For frequentist lasso regression (Tibshirani 1996) the latter penalty produces sparse regression fits, in that the estimated coefficients become exactly zero. In the Bayesian case the Bayes estimates do not provide exact annihilation of coefficients, but produce approximate sparseness (Park and Casella 2008).

At this point we note that, during the years following the appearance of Park and Casella (2008), several other proposals for sparse shrinkage of the β_j appeared in the literature (e.g., Armagan, Dunson, and Lee 2013; Carvalho, Polson, and Scott 2010; Griffin and Brown 2011). Their accommodation in *Infer.NET* could also be entertained. Here we restrict attention to Laplace-based shrinkage.

We commence with the goal of fitting the following model in *Infer.NET*:

$$\begin{aligned} \mathbf{y} | \beta_0, \beta, \tau_\epsilon &\sim N(\mathbf{1}\beta_0 + \mathbf{X}\beta, \tau_\epsilon^{-1}\mathbf{I}), \quad \beta_j | \tau_\beta \stackrel{\text{ind.}}{\sim} \text{Laplace}(0, \tau_\beta^{-1}) \\ \beta_0 &\sim N(0, \tau_{\beta_0}^{-1}), \quad \tau_\beta \sim \text{Half-Cauchy}(A_\beta), \quad \tau_\epsilon \sim \text{Half-Cauchy}(A_\epsilon). \end{aligned} \quad (43)$$

The current release of *Infer.NET* does not support the Laplace distribution, so we require an auxiliary variable representation in a similar vein to what was used for the t distribution in Section 5. We first note the distributional statement

$$x | \tau \sim \text{Laplace}(0, \tau^{-1}) \quad \text{if and only if} \quad x | g \sim N(0, g), \quad g | \tau \sim \text{Gamma}(1, \frac{1}{2} \tau), \quad (44)$$

which is exploited by Park and Casella (2008). However, auxiliary variable representation (44) is also not supported by *Infer.NET*, due to violation of its conjugacy rules. We get around this by working with the alternative auxiliary variable set-up:

$$x | c \sim N(0, 1/c), \quad c | d \sim \text{Gamma}(M, M d), \quad d | \tau \sim \text{Gamma}(1, \frac{1}{2} \tau), \quad (45)$$

which is supported by *Infer.NET*, and leads to good approximation to the $\text{Laplace}(0, \tau^{-1})$ distribution when M is large. For any given $M > 0$, the density function of x satisfying auxiliary variable representation (45) is

$$\begin{aligned} p(x | \tau; M) &= \int_0^\infty \int_0^\infty (2\pi/c)^{-1/2} \exp\left(-\frac{1}{2} c x\right) \frac{(M d)^M}{\Gamma(M)} c^{M-1} \exp(-M d c) \\ &\quad \times \frac{1}{2} \tau \exp\left(-\frac{1}{2} d \tau\right) dc dd \\ &= \frac{\tau^{1/2} \Gamma\left(M + \frac{1}{2}\right)}{2\Gamma(M) M^{1/2}} {}_1F_1\left(M + \frac{1}{2}; \frac{1}{2}; \frac{\tau x^2}{4M}\right) - \frac{1}{2} \tau^{1/2} |x| {}_1F_1\left(M + 1; 3/2; \frac{\tau x^2}{4M}\right) \end{aligned}$$

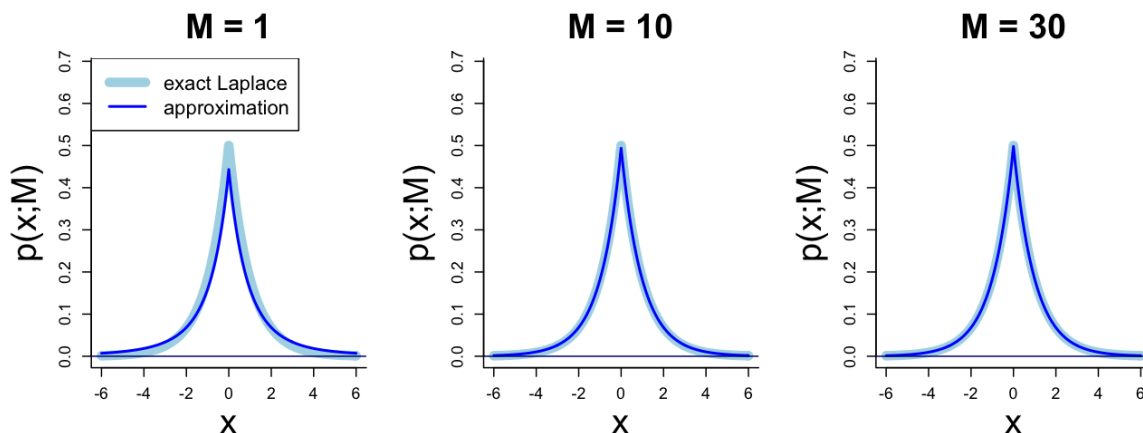


Figure 10: The exact Laplace(0, 1) density function $p(x) = \frac{1}{2} \exp(-|x|)$ and three approximations to it based on $p(x|\tau; M)$ for $\tau = 1$ and $M = 1, 10, 30$.

where ${}_1F_1$ denotes the degenerate hypergeometric function (e.g., [Gradshteyn and Ryzhik 1994](#)). Figure 10 shows $p(x|\tau; M)$ for $\tau = 1$ and $M = 1, 10, 30$. The approximation to the Laplace(0, 1) density function is excellent for M as low as 10.

The convergence in distribution of random variables having the density function $p(x|\tau; M)$ to a Laplace(0, τ^{-1}) random variable as $M \rightarrow \infty$ may be established using (28) and the fact that the family of t distributions tends to a Normal distribution as the degrees of freedom parameter increases, and then mixing this limiting distribution with a Gamma($1, \frac{1}{2} \tau$) variable. In lieu of (43), the actual model fitted in Infer.NET is:

$$\begin{aligned}
 y|\beta_0, \beta, \tau_\varepsilon &\sim N(1\beta_0 + X\beta, \tau_\varepsilon^{-1}I), \\
 \mathbf{a}|\beta_0, \beta, \mathbf{c}_1, \dots, \mathbf{c}_p &\sim N\left(\begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}, \begin{bmatrix} \tau_{\beta_0}^{-1} & \mathbf{0} \\ \mathbf{0} & \text{diag}(\mathbf{1}/\mathbf{c}_j)I \end{bmatrix}\right), \\
 \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} &\sim N(0, \kappa^{-1}I), \quad c_j|d_j \stackrel{\text{ind.}}{\sim} \text{Gamma}(M, M d_j), \tag{46}
 \end{aligned}$$

$$d_j|\tau_\beta \stackrel{\text{ind.}}{\sim} \text{Gamma}(1, \frac{1}{2}\tau_\beta), \quad \tau_\beta|b_\beta \sim \text{Gamma}(\frac{1}{2}, b_\beta), \quad b_\beta \sim \text{Gamma}(\frac{1}{2}, 1/A_\beta^2),$$

$$\tau_\varepsilon|b_\varepsilon \sim \text{Gamma}(\frac{1}{2}, b_\varepsilon), \quad b_\varepsilon \sim \text{Gamma}(\frac{1}{2}, 1/A_\varepsilon^2).$$

We use $\kappa = 10^{-10}$, $M = 100$, $\tau_{\beta_0} = 10^{-10}$, $A_\beta = A_\varepsilon = 10^5$.

Our illustration of (46) involves data from a diabetes study that was also used in [Park and Casella \(2008\)](#). The sample size is $n = 442$ and the number of predictor variables is $p = 10$. The response variable is a continuous index of disease progression one year after baseline. A description of the predictor variables is given in [Efron, Hastie, Johnstone, and Tibshirani \(2004\)](#). Their abbreviations, used in [Park and Casella \(2008\)](#), are: (1) age, (2) sex, (3) bmi, (4) map, (5) tc, (6) ldl, (7) hdl, (8) tch, (9) ltg, and (10) glu. The number of mean field variational Bayes iterations was fixed at 100. The number of burn-in samples and samples used for inference by the MCMC fit were both set to 5000.

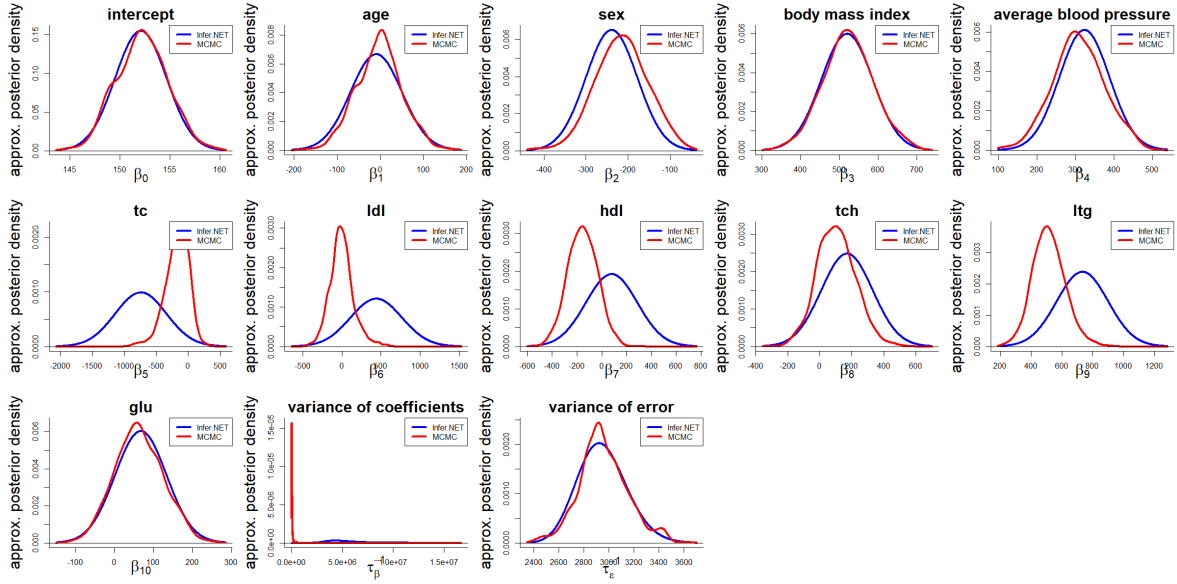


Figure 11: Approximate posterior density functions of model parameters in fitting (46) to data from a diabetes study.

Figure 11 shows the fits obtained from both Infer.NET and MCMC. We see that the correspondence is very good, especially for the regression coefficients.

9. Measurement error model

In some situations some variables may be measured with error, i.e., the observed data are not the true values of those variables themselves, but a contaminated version of these variables. Examples of such situations include AIDS studies where CD4 counts are known to be measured with errors (Wu 2002) or air pollution data (Bachrach *et al.* 1999). Carroll, Ruppert, Stefanski, and Crainiceanu (2006) offer a comprehensive treatment of measurement error models. As described there, failure to account for measurement error can result in biased and misleading conclusions.

Let (x_i, y_i) , $1 \leq i \leq n$, be a set of predictor/response pairs that are modeled according to

$$y_i | x_i, \beta_0, \beta_1, \tau_\varepsilon \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \tau_\varepsilon^{-1}), \quad 1 \leq i \leq n. \quad (47)$$

However, instead of observing the x_i s we observe

$$w_i | x_i \stackrel{\text{ind.}}{\sim} N(x_i, \tau_z^{-1}), \quad 1 \leq i \leq n. \quad (48)$$

In other words, the predictors are measured with error with τ_z controlling the extent of the contamination. In general τ_z can be estimated through validation data, where some of the x_i values are observed, or replication data, where replicates of the w_i values are available. To simplify exposition we will assume that τ_z is known. Furthermore, we will assume that the x_i s are Normally distributed:

$$x_i | \mu_x, \tau_x \stackrel{\text{ind.}}{\sim} N(\mu_x, \tau_x^{-1}), \quad 1 \leq i \leq n, \quad (49)$$

where μ_x and τ_x are additional parameters to be estimated.

Combining (47), (48) and (49) and, including aforementioned priors for variances, a Bayesian measurement error model can be represented by:

$$\begin{aligned} \mathbf{y} | \mathbf{x}, \boldsymbol{\beta}, \tau_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta}, \tau_\varepsilon^{-1}\mathbf{I}), & \mathbf{w} | \mathbf{x} &\sim N(\mathbf{x}, \tau_z^{-1}\mathbf{I}), & \mathbf{x} | \mu_x, \tau_x &\sim N(\mu_x \mathbf{1}, \tau_x^{-1}\mathbf{I}), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), & \tau_\varepsilon^{-1/2} &\sim \text{Half-Cauchy}(\mathbf{A}_\varepsilon), \\ \mu_x &\sim N(0, \tau_{\mu_x}^{-1}), & \text{and } \tau_x^{-1/2} &\sim \text{Half-Cauchy}(A_x), \end{aligned} \quad (50)$$

where $\boldsymbol{\beta} = [\beta_0, \beta_1]$, $\mathbf{X} = [\mathbf{1}, \mathbf{x}]$, τ_β , τ_μ , A_ε and A_x are user-specified constants and the vectors \mathbf{y} and \mathbf{w} are observed. We set $\tau_\beta = \tau_\mu = 10^{-10}$ and $A_\varepsilon = A_x = 10^5$.

Invoking (1), we arrive at the model

$$\begin{aligned} \mathbf{y} | \mathbf{x}, \boldsymbol{\beta}, \tau_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta}, \tau_\varepsilon^{-1}\mathbf{I}), & \mathbf{w} | \mathbf{x} &\sim N(\mathbf{x}, \tau_z^{-1}\mathbf{I}), & \mathbf{x} | \mu_x, \tau_x &\sim N(\mu_x \mathbf{1}, \tau_x^{-1}\mathbf{I}) \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \tau_\beta^{-1}\mathbf{I}), & \tau_\varepsilon | b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, 1/b_\varepsilon), & b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2), \\ \mu_x &\sim N(0, \tau_{\mu_x}^{-1}), & \tau_x | b_x &\sim \text{Gamma}(\tfrac{1}{2}, 1/b_x) & \text{and } b_x &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_x^2). \end{aligned} \quad (51)$$

Infer.NET 2.6 does not appear to be able to fit the model (51) under the product restriction

$$\mathbf{q}(\boldsymbol{\beta}, \tau_\varepsilon, \mathbf{b}_\varepsilon, \mathbf{x}, \mu_x, \tau_x, \mathbf{b}_x) = \mathbf{q}(\beta_0, \beta_1) \mathbf{q}(\mu_x) \mathbf{q}(\mathbf{x}) \mathbf{q}(\tau_\varepsilon, \tau_x) \mathbf{q}(\mathbf{b}_\varepsilon, \mathbf{b}_x).$$

However, Infer.NET was able to fit this model with

$$\mathbf{q}(\boldsymbol{\beta}, \tau_\varepsilon, \mathbf{b}_\varepsilon, \mathbf{x}, \mu_x, \tau_x, \mathbf{b}_x) = \mathbf{q}(\beta_0) \mathbf{q}(\beta_1) \mathbf{q}(\mu_x) \mathbf{q}(\mathbf{x}) \mathbf{q}(\tau_\varepsilon, \tau_x) \mathbf{q}(\mathbf{b}_\varepsilon, \mathbf{b}_x). \quad (52)$$

Unfortunately, fitting under restriction (52) leads to poor accuracy. One possible remedy involves the centering transformation:

$$\tilde{w}_i = w_i - \bar{w}.$$

Under this transformation, and with diffuse priors, the marginal posterior distributions of β_0 and β_1 are nearly independent.

Let $\tilde{q}(\beta_0)$, $\tilde{q}(\beta_1)$, $\tilde{q}(\mu_x)$, $\tilde{q}(\mathbf{x})$, $\tilde{q}(\tau_\varepsilon, \tau_x)$ and $\tilde{q}(b_\varepsilon, b_x)$ be the optimal values of $q(\beta_0)$, $q(\beta_1)$, $q(\mu_x)$, $q(\mathbf{x})$, $q(\tau_\varepsilon, \tau_x)$ and $q(b_\varepsilon, b_x)$ when the transformed \tilde{w}_i s are used in place of the w_i s. This corresponds to the transformation on the \mathbf{X} matrix

$$\tilde{\mathbf{X}} = \mathbf{X}\mathbf{R}, \quad \text{where} \quad \mathbf{R} = \begin{bmatrix} 1 & 0 \\ -\bar{w} & 1 \end{bmatrix}.$$

Suppose that

$$\begin{aligned} \tilde{q}(\beta_0) &\sim N(\tilde{\mu}_{q(\beta_0)}, \tilde{\sigma}_{q(\beta_0)}^2), & \tilde{q}(\beta_1) &\sim N(\tilde{\mu}_{q(\beta_1)}, \tilde{\sigma}_{q(\beta_1)}^2), \\ \tilde{q}(\mu_x) &\sim N(\tilde{\mu}_{q(\mu_x)}, \tilde{\sigma}_{q(\mu_x)}^2), & \text{and } \tilde{q}(x_i) &\sim N(\tilde{\mu}_{q(x_i)}, \tilde{\sigma}_{q(x_i)}^2), \quad 1 \leq i \leq n. \end{aligned}$$

Then we can back-transform approximately using

$$\begin{aligned} q(\beta_0, \beta_1) &\sim N(\boldsymbol{\mu}_{q(\boldsymbol{\beta})}, \boldsymbol{\Sigma}_{q(\boldsymbol{\beta})}), & q(\mu_x) &\sim N(\mu_{q(\mu_x)}, \sigma_{q(\mu_x)}^2), \\ q(x_i) &\sim N(\mu_{q(x_i)}, \sigma_{q(x_i)}^2), & 1 \leq i &\leq n, \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\beta)} &= R \begin{bmatrix} \tilde{\mu}_{q(\beta_0)} \\ \tilde{\mu}_{q(\beta_1)} \end{bmatrix}, & \boldsymbol{\Sigma}_{q(\beta)} &= R \begin{bmatrix} \tilde{\sigma}_{q(\beta_0)}^2 & 0 \\ 0 & \tilde{\sigma}_{q(\beta_1)}^2 \end{bmatrix} R^\top, \\ \mu_{q(\mu_x)} &= \tilde{\mu}_{q(\mu_x)} + \bar{w}, & \sigma_{q(\mu_x)}^2 &= \tilde{\sigma}_{q(\mu_x)}^2, \\ \mu_{q(x_i)} &= \tilde{\mu}_{q(x_i)} + \bar{w}, & \sigma_{q(x_i)}^2 &= \tilde{\sigma}_{q(x_i)}^2, \quad 1 \leq i \leq n, \end{aligned}$$

$$q(\tau_\varepsilon, \tau_x) = \tilde{q}(\tau_\varepsilon, \tau_x) \text{ and } q(b_\varepsilon, b_x) = \tilde{q}(b_\varepsilon, b_x).$$

To illustrate the use of *Infer.NET* we simulated data according to

$$x_i \stackrel{\text{ind.}}{\sim} N(1/2, 1/36), \quad w_i | x_i \stackrel{\text{ind.}}{\sim} N(x_i, 1/100) \quad \text{and} \quad y_i | x_i \stackrel{\text{ind.}}{\sim} N(0.3 + 0.7x_i, 1/16) \quad (53)$$

for $1 \leq i \leq n$ with $n = 50$. Results from a single simulation are illustrated in Figure 12. Similar results were obtained from other simulated data sets. From Figure 12 we see reasonable agreement of posterior density estimates produced by *Infer.NET* and MCMC for all parameters. The number of mean field variational Bayes iterations was fixed at 100. The number of burn-in samples was and the samples used for inference by the MCMC fit was set to 50000. Extension to the case where the mean of the y_i s are modeled nonparametrically is covered in [Pham, Ormerod, and Wand \(2013\)](#). However, the current version of *Infer.NET* does not support models of this type. Such versatility limitations of *Infer.NET* are discussed in Section 11.1.

10. Missing predictor model

Missing data is a commonly occurring issue in statistical problems. Naïve methods for dealing with this issue, such as ignoring samples containing missing values, can be shown to be demonstrably optimistic (i.e., standard errors are deflated), biased or simply an inefficient use of the data. [Little and Rubin \(2002\)](#) contains a detailed discussion of the various issues at stake.

Consider a simple linear regression model with response y_i with predictor x_i :

$$y_i | x_i, \beta_0, \beta_1, \tau_\varepsilon \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \tau_\varepsilon^{-1}), \quad 1 \leq i \leq n. \quad (54)$$

Suppose that some of the x_i s are missing and let r_i be an indicator for x_i being observed. Specifically,

$$r_i = \begin{cases} 1, & \text{if } x_i \text{ is observed,} \\ 0, & \text{if } x_i \text{ is missing,} \end{cases} \quad \text{for } 1 \leq i \leq n.$$

A missing data model contains at least two components:

- A model for the underlying distribution of the variable that is subject to missing data. We will use

$$x_i | \mu_x, \tau_x \stackrel{\text{ind.}}{\sim} N(\mu_x, \tau_x^{-1}), \quad 1 \leq i \leq n. \quad (55)$$

- A model for the missing data mechanism, which models the probability distribution of r_i .

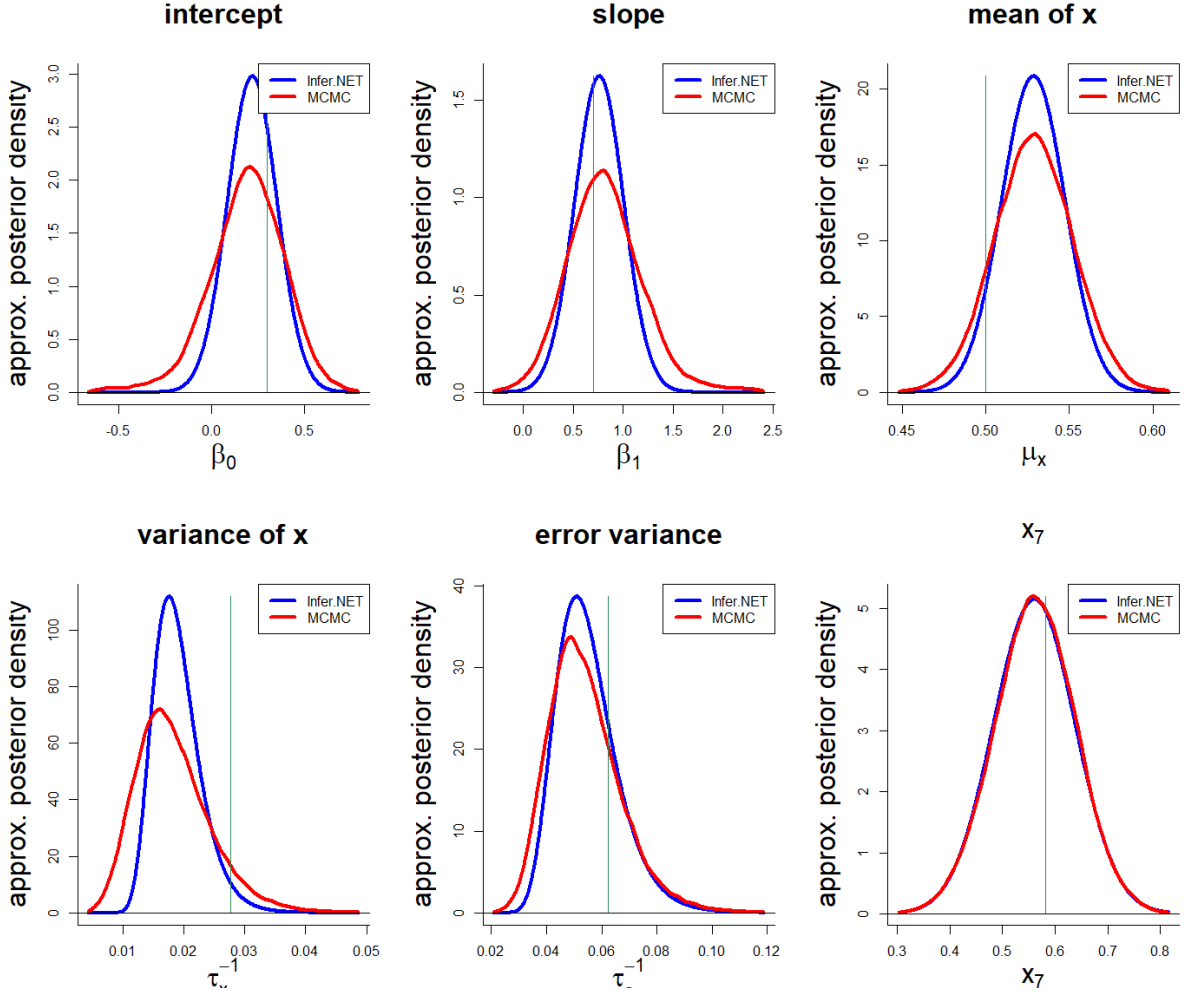


Figure 12: Variational Bayes approximate posterior density functions produced by Infer.NET for simulated data (53).

As detailed in [Faes et al. \(2011\)](#), there are several possible models, with varying degrees of sophistication, for the missing data mechanism. Here the simplest such model, known as *missing completely at random*, is used:

$$\begin{aligned}
 \mathbf{y}|\mathbf{x}, \boldsymbol{\beta}, \boldsymbol{\tau}_\varepsilon &\sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\tau}_\varepsilon^{-1}\mathbf{I}), & \mathbf{x}|\boldsymbol{\mu}_x, \tau_x &\sim N(\boldsymbol{\mu}_x\mathbf{1}, \tau_x^{-1}\mathbf{I}), & r_i &\stackrel{\text{ind.}}{\sim} \text{Bernoulli}(p), \\
 \boldsymbol{\beta} &\sim N(\mathbf{0}, \boldsymbol{\tau}_\beta^{-1}\mathbf{I}), & \tau_\varepsilon|b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, 1/b_\varepsilon), & b_\varepsilon &\sim \text{Gamma}(\tfrac{1}{2}, 1/A_\varepsilon^2), \\
 \mu_x &\sim N(0, \tau_\mu^{-1}), & \tau_x|b_x &\sim \text{Gamma}(\tfrac{1}{2}, b_x) & \text{and} & b_x \sim \text{Gamma}(\tfrac{1}{2}, 1/A_x^2),
 \end{aligned} \tag{56}$$

where $\boldsymbol{\beta} = [\beta_0, \beta_1]$, $\mathbf{X} = [\mathbf{1}, \mathbf{x}]$, τ_β , τ_μ , A_ε and A_x are user-specified hyperparameters. Again, we will use $\tau_\beta = \tau_\mu = 10^{-10}$ and $A_\varepsilon = A_x = 10^5$.

The simulated data that we use in our illustration of Infer.NET involves the following sample size and ‘true values’:

$$n = 50, \quad \beta_0 = 0.3, \quad \beta_1 = 0.7, \quad \mu_x = 0.5, \quad \tau_x = 36, \quad \text{and} \quad p = 0.75. \tag{57}$$

Putting $p = 0.75$ implies that, on average, 25% of the predictor values are missing completely at random.

Infer.NET 2.6 is not able to fit the model (51) under the product restriction

$$q(\beta, \tau_\varepsilon, \mathbf{b}_\varepsilon, \mathbf{x}_{\text{mis}}, \boldsymbol{\mu}_x, \boldsymbol{\tau}_x, \mathbf{b}_x) = q(\beta_0, \beta_1) q(\mathbf{x}_{\text{mis}}) q(\tau_\varepsilon, \tau_x) q(\mathbf{b}_\varepsilon, \mathbf{b}_x), \quad (58)$$

where \mathbf{x}_{mis} denotes the vector of missing predictors, but is able to handle

$$q(\beta, \tau_\varepsilon, \mathbf{b}_\varepsilon, \mathbf{x}_{\text{mis}}, \boldsymbol{\mu}_x, \boldsymbol{\tau}_x, \mathbf{b}_x) = q(\beta_0) q(\beta_1) q(\mathbf{x}_{\text{mis}}) q(\tau_\varepsilon, \tau_x) q(\mathbf{b}_\varepsilon, \mathbf{b}_x), \quad (59)$$

and was used in our illustration of *Infer.NET* for such models.

We also used BUGS to perform Bayesian inference via MCMC as a benchmark for comparison of results. For similar reasons as in the previous section fitting under this restriction (59) leads to poor results. Instead we perform the centering transformation

$$\tilde{x}_i = x_i - \bar{x}_{\text{obs}},$$

where \bar{x}_{obs} is the mean of the observed x values. This transformation makes the marginal posterior distributions of β_0 and β_1 nearly independent.

Let $\tilde{q}(\beta_0)$, $\tilde{q}(\beta_1)$, $\tilde{q}(\mu_x)$, $\tilde{q}(x_{\text{mis}})$, $\tilde{q}(\tau_\varepsilon, \tau_x)$ and $\tilde{q}(b_\varepsilon, b_x)$ be the optimal values of $q(\beta_0)$, $q(\beta_1)$, $q(\mu_x)$, $q(x_{\text{mis}})$, $q(\tau_\varepsilon, \tau_x)$ and $q(b_\varepsilon, b_x)$ when the transformed \tilde{x}_i s are used in place of the x_i s. This corresponds to the following transformation of the X matrix:

$$\tilde{X} = XR, \quad \text{where} \quad R = \begin{bmatrix} 1 & 0 \\ -\bar{x}_{\text{obs}} & 1 \end{bmatrix}.$$

Suppose that

$$\begin{aligned} \tilde{q}(\beta_0) &\sim N(\tilde{\mu}_{q(\beta_0)}, \tilde{\sigma}_{q(\beta_0)}^2), & \tilde{q}(\beta_1) &\sim N(\tilde{\mu}_{q(\beta_1)}, \tilde{\sigma}_{q(\beta_1)}^2), \\ \tilde{q}(\mu_x) &\sim N(\tilde{\mu}_{q(\mu_x)}, \tilde{\sigma}_{q(\mu_x)}^2), & \text{and} & \quad \tilde{q}(x_i) \sim N(\tilde{\mu}_{q(x_{\text{mis},i})}, \tilde{\sigma}_{q(x_{\text{mis},i})}^2), \quad 1 \leq i \leq n_{\text{mis}}, \end{aligned}$$

which we back-transform approximately using

$$\begin{aligned} q(\beta_0, \beta_1) &\sim N(\boldsymbol{\mu}_{q(\beta)}, \boldsymbol{\Sigma}_{q(\beta)}), & q(\mu_x) &\sim N(\mu_{q(\mu_x)}, \sigma_{q(\mu_x)}^2), \\ q(x_{\text{mis},i}) &\sim N(\mu_{q(x_{\text{mis},i})}, \sigma_{q(x_{\text{mis},i})}^2), \quad 1 \leq i \leq n, \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\mu}_{q(\beta)} &= R \begin{bmatrix} \tilde{\mu}_{q(\beta_0)} \\ \tilde{\mu}_{q(\beta_1)} \end{bmatrix}, & \boldsymbol{\Sigma}_{q(\beta)} &= R \begin{bmatrix} \tilde{\sigma}_{q(\beta_0)}^2 & 0 \\ 0 & \tilde{\sigma}_{q(\beta_1)}^2 \end{bmatrix} R^\top, \\ \mu_{q(\mu_x)} &= \tilde{\mu}_{q(\mu_x)} + \bar{x}_{\text{obs}}, & \sigma_{q(\mu_x)}^2 &= \tilde{\sigma}_{q(\mu_x)}^2, \\ \mu_{q(x_{\text{mis},i})} &= \tilde{\mu}_{q(x_{\text{mis},i})} + \bar{x}_{\text{obs}}, & \sigma_{q(x_{\text{mis},i})}^2 &= \tilde{\sigma}_{q(x_{\text{mis},i})}^2, \quad 1 \leq i \leq n, \end{aligned}$$

$$q(\tau_\varepsilon, \tau_x) = \tilde{q}(\tau_\varepsilon, \tau_x) \text{ and } q(b_\varepsilon, b_x) = \tilde{q}(b_\varepsilon, b_x).$$

Results from a single simulation are illustrated in Figure 13. Here the number of mean field variational Bayes iterations was fixed at 100, and for MCMC fits the number of burn-in samples was 5000 and the number of samples used for inference was set to 50000. Similar results

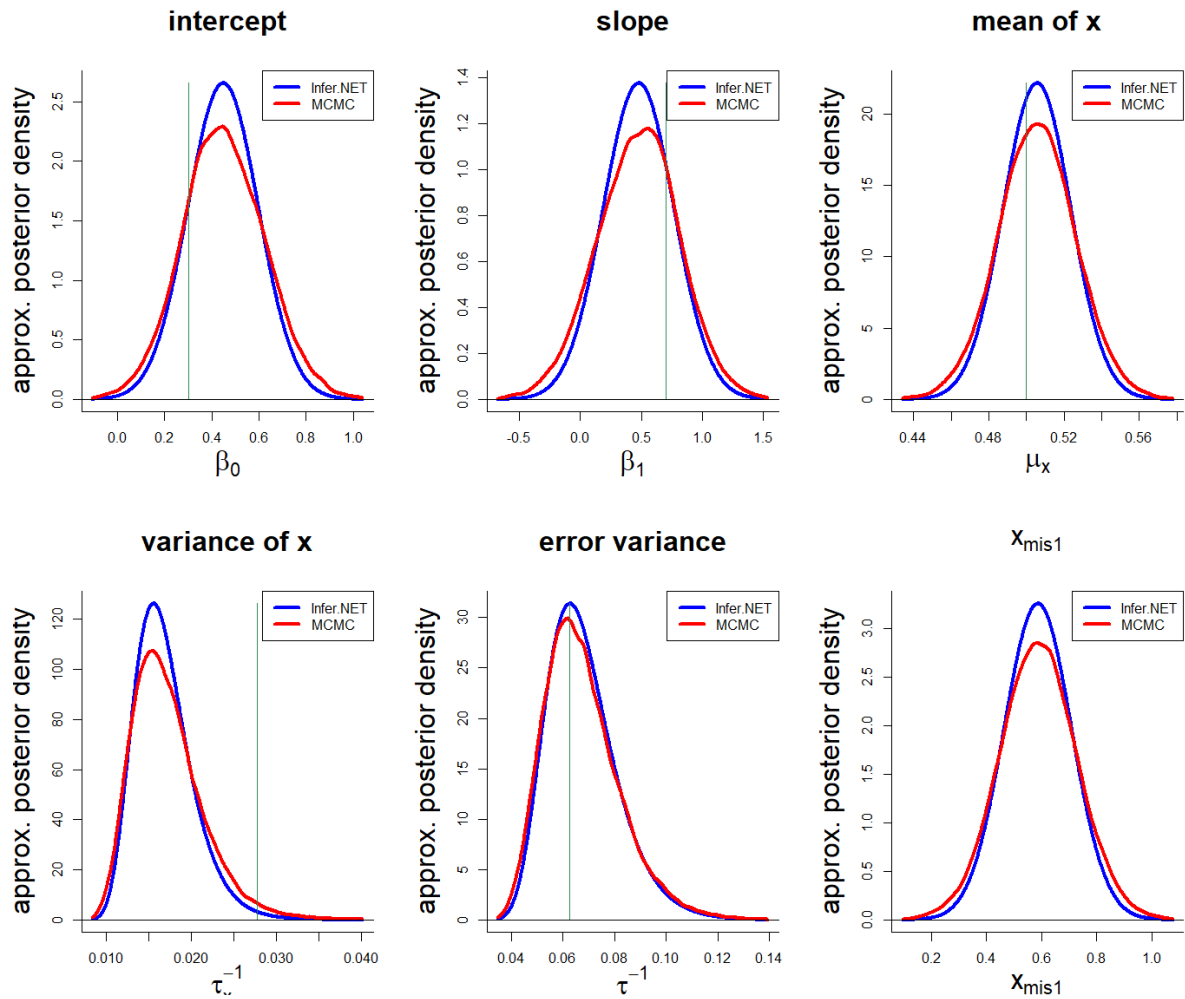


Figure 13: Variational Bayes approximate posterior density functions produced by Infer.NET data simulated according to (57).

were obtained from other simulated data sets. From Figure 13 we see reasonable agreement of posterior density estimates produced by Infer.NET and MCMC for all parameters. Extension to the case where the mean of y is modeled nonparametrically is covered in Faes *et al.* (2011). However, the current version of Infer.NET does not support models of this type.

11. Comparison with BUGS

We now make some comparisons between Infer.NET and BUGS in the context of Bayesian semiparametric regression.

11.1. Versatility comparison

These comments mainly echo those given in Section 5 of Wang and Wand (2011), and thus are quite brief. BUGS is much more versatile than Infer.NET, with the latter being subject

Section number	Semiparametric regression model	Time in seconds for <i>Infer.NET</i>	Time in seconds for <i>BUGS</i>
3	Simple semiparametric regression	3.14 (0.10)	6.98 (0.01)
4	Generalized additive model (logistic)	4.22 (0.04)	122.48 (5.91)
4	Generalized additive model (probit)	8.45 (0.07)	81.13 (3.63)
5	Robust nonparametric regression	103.61 (2.20)	22.15 (0.10)
6	Semiparametric mixed model	594.99 (3.60)	323.62 (0.70)
7	Geoadditive model	1383.63 (4.25)	— —
8	Bayesian lasso regression	3.15 (0.05)	253.19 (1.44)
9	Measurement error model	3.72 (0.06)	5.51 (0.16)
10	Missing predictor model	5.57 (0.08)	6.11 (0.09)

Table 3: Average run times (standard errors) in seconds over 20 runs of the methods for each of the examples in the paper.

to restrictions such as standard distributional forms, conjugacy rules and not being able to handle models such as (14) without the trick manifest in (16). The measurement error and missing data regression examples given in Sections 9 and 10 are feasible in *Infer.NET* for parametric regression, but not for semiparametric regression such as the examples in Sections 4, 5 and 8 of Marley and Wand (2010). Nevertheless, as demonstrated in this article, there is a wide variety of semiparametric regression models that can be handled by *Infer.NET*.

11.2. Accuracy comparison

In general, *BUGS* can always be more accurate than *Infer.NET* since MCMC suffers only from Monte Carlo error, which can be made arbitrarily small via larger sample sizes. The *Infer.NET* inference engines, expectation propagation and variational message passing, have inherent approximation errors that cannot be completely eliminated. However, our examples show that the accuracy of *Infer.NET* is quite good for a wide variety of semiparametric regression models.

11.3. Timing comparison

Table 3 gives some indication of the relative computing times for the examples in the paper. In this table we report the average elapsed (and standard error) of the computing times over 100 runs with the number of variational Bayes or expectation propagation iterations set to 100 and the MCMC sample sizes set to 10000. This was sufficient for convergence in these particular examples. All examples were run on the third author’s laptop computer (64 bit Windows 10 Intel i7-7600MX central processing unit at 2.8GHz with 2 hyperthreaded cores and 32GB of random access memory). We concede that comparison of deterministic and Monte Carlo algorithms is fraught with difficulties. However, convergence criteria aside, these times give an indication of the performance in terms of the implementations of each of the algorithms for particular models on a fast 2018 computer. Note that we did not fit the geoadditive model in Section 7 via *BUGS* due to the excessive time required.

Table 3 reveals that *Infer.NET* offers considerable speed-ups compared with *BUGS* for most of the examples. The exceptions are the robust nonparametric regression model of Section 5

and the semiparametric mixed model of Section 6. The slowness of the robust nonparametric regression fit is mainly explained by the multiple calls to Infer.NET, corresponding to the degrees of freedom grid. The semiparametric mixed model is quite slow in the current release of Infer.NET due to the full sparse design matrices being carried around in the calculations. Recently developed streamlined approaches to variational inference for longitudinal and multilevel data analysis (Lee and Wand 2016) offer significant speed-ups. The geoaddivitive model of Section 7 is also slow to fit, with the large design matrices being a likely reason.

12. Summary

Through several examples, the efficacy of Infer.NET for semiparametric regression analysis has been demonstrated. Generally speaking, the fitting and inference is shown to be quite fast and accurate. Models with very large design matrices are an exception and can take significant amounts of time to fit using the current release of Infer.NET.

Our survey of Infer.NET in the context of semiparametric regression will aid future analyses of this type via fast Bayesian inference software, by building upon the examples presented here. It may also influence future directions for the development of fast Bayesian methodology and software.

Acknowledgments

This research was partially funded by the Australian Research Council Discovery Project DP110100061. We are grateful to Andrew Chernih for his provision of the Sydney property data.

References

- Armagan A, Dunson DB, Lee J (2013). “Generalized Double Pareto Shrinkage.” *Statistica Sinica*, **23**(1), 119–143. doi:10.5705/ss.2011.048.
- Bachrach LK, Hastie T, Wang MC, Narasimhan B, Marcus R (1999). “Bone Mineral Acquisition in Healthy Asian, Hispanic, Black, and Caucasian Youth: A Longitudinal Study.” *Journal of Clinical Endocrinology & Metabolism*, **84**(12), 4702–4712. doi:10.1210/jc.84.12.4702.
- Carroll RJ, Ruppert D, Stefanski LA, Crainiceanu C (2006). *Measurement Error in Nonlinear Models: A Modern Perspective*. 2nd edition. Chapman and Hall, London. doi:10.1201/9781420010138.
- Carvalho CM, Polson NG, Scott JG (2010). “The Horseshoe Estimator for Sparse Signals.” *Biometrika*, **97**(2), 465–480. doi:10.1093/biomet/asq017.
- Efron B, Hastie T, Johnstone I, Tibshirani R (2004). “Least Angle Regression.” *The Annals of Statistics*, **32**(2), 407–451. doi:10.1214/009053604000000067.
- Eilers PHC, Marx BD (1996). “Flexible Smoothing with B -Splines and Penalties.” *Statistical Science*, **11**(2), 89–121. doi:10.1214/ss/1038425655.

- Faes C, Ormerod JT, Wand MP (2011). “Variational Bayesian Inference for Parametric and Nonparametric Regression with Missing Data.” *Journal of the American Statistical Association*, **106**(495), 959–971. doi:10.1198/jasa.2011.tm10301.
- Fasshauer GE (2007). *Meshfree Approximation Methods with MATLAB*, volume 6 of *Interdisciplinary Mathematical Sciences*. World Scientific, Singapore.
- Gelman A (2006). “Prior Distributions for Variance Parameters in Hierarchical Models.” *Bayesian Analysis*, **1**(3), 515–534. doi:10.1214/06-ba117a.
- Gradshteyn IS, Ryzhik IM (1994). *Table of Integrals, Series, and Products*. Academic Press, Boston.
- Griffin JE, Brown PJ (2011). “Bayesian Hyper-Lassos with Non-Convex Penalization.” *Australian and New Zealand Journal of Statistics*, **53**(4), 423–442. doi:10.1111/j.1467-842x.2011.00641.x.
- Haberman SJ (1976). “Generalized Residuals for Log-Linear Models.” In *Proceedings of the 9th International Biometrics Conference*, pp. 104–122. Boston.
- Kammann EE, Wand MP (2003). “Geoadditive Models.” *Journal of the Royal Statistical Society C*, **52**(1), 1–18. doi:10.1111/1467-9876.00385.
- Kaufman L, Rousseeuw PJ (1990). *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley & Sons, New York. doi:10.1002/9780470316801.
- Kim ASI, Wand MP (2016). “The Explicit Form of Expectation Propagation for a Simple Statistical Model.” *Electronic Journal of Statistics*, **10**(1), 550–581. doi:10.1214/16-ejs1114.
- Knowles DA, Minka TP (2011). “Non-Conjugate Message Passing for Multinomial and Binary Regression.” In J Shawe-Taylor, RS Zemel, P Bartlett, F Pereira, KQ Weinberger (eds.), *Advances in Neural Information Processing Systems 24*, pp. 1701–1709. Neural Information Processing System Foundation, Inc.
- Lange KL, Little RJA, Taylor JMG (1989). “Robust Statistical Modeling Using the t -Distribution.” *The Journal of the American Statistical Association*, **84**(408), 881–896. doi:10.1080/01621459.1989.10478852.
- Lee YY, Wand MP (2016). “Streamlined Mean Field Variational Bayes For Longitudinal and Multilevel Data Analysis.” *Biometrical Journal*, **58**(4), 868–895. doi:10.1002/bimj.201500007.
- Little RJA, Rubin DB (2002). *Statistical Analysis with Missing Data*. 2nd edition. John Wiley & Sons, New York. doi:10.1002/9781119013563.
- Lunn D, Thomas A, Best NG, Spiegelhalter DJ (2000). “**WinBUGS** – A Bayesian Modelling Framework: Concepts, Structure, and Extensibility.” *Statistics and Computing*, **10**(4), 325–337. doi:10.1023/a:1008929526011.
- Marley JK, Wand MP (2010). “Non-Standard Semiparametric Regression via **BRugs**.” *Journal of Statistical Software*, **37**(5), 1–30. doi:10.18637/jss.v037.i05.

- Minka T (2001). “Expectation Propagation for Approximate Bayesian Inference.” *Proceedings of Conference on Uncertainty in Artificial Intelligence*, pp. 51–76.
- Minka T (2005). “Divergence Measures and Message Passing.” *Technical Report MSR-TR-2008-173*, Microsoft Research Technical Report Series.
- Minka T, Winn J (2008). “Gates: A Graphical Notation for Mixture Models.” *Technical Report MSR-TR-2008-185*, Microsoft Research Technical Report Series.
- Minka T, Winn JM, Guiver JP, Knowles D (2014). *Infer.NET 2.6*. Microsoft Research Cambridge, URL <http://research.microsoft.com/infernet>.
- Ormerod JT, Wand MP (2010). “Explaining Variational Approximations.” *The American Statistician*, **64**(2), 140–153. doi:10.1198/tast.2010.09058.
- Park T, Casella G (2008). “The Bayesian Lasso.” *Journal of the American Statistical Association*, **103**(482), 681–686. doi:10.1198/016214508000000337.
- Pham T, Ormerod JT, Wand MP (2013). “Mean Field Variational Bayesian Inference for Nonparametric Regression with Measurement Error.” *Computational Statistics & Data Analysis*, **68**, 375–387. doi:10.1016/j.csda.2013.07.014.
- R Core Team (2018). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Ruppert D, Wand MP, Carroll RJ (2003). *Semiparametric Regression*. Cambridge University Press, New York.
- Ruppert D, Wand MP, Carroll RJ (2009). “Semiparametric Regression During 2003–2007.” *Electronic Journal of Statistics*, **3**, 1193–1265. doi:10.1214/09-ejs525.
- Saul LK, Jordan MI (1996). “Exploiting Tractable Substructures in Intractable Networks.” In *Advances in Neural Information Processing Systems*, pp. 435–442. MIT Press, Cambridge, Massachusetts.
- Staudenmayer J, Lake EE, Wand MP (2009). “Robustness for General Design Mixed Models Using the t -Distribution.” *Statistical Modelling*, **9**(3), 235–255. doi:10.1177/1471082x0800900304.
- Tibshirani R (1996). “Regression Shrinkage and Selection via the Lasso.” *Journal of the Royal Statistical Society B*, **58**(1), 267–288.
- Wainwright M, Jordan MI (2008). “Graphical Models, Exponential Families, and Variational Inference.” *Foundations and Trends in Machine Learning*, **1**(1–2), 1–305. doi:10.1561/2200000001.
- Wand MP (2009). “Semiparametric Regression and Graphical Models.” *Australian and New Zealand Journal of Statistics*, **51**(1), 9–41. doi:10.1111/j.1467-842x.2009.00538.x.
- Wand MP (2015). **KernSmooth**: *Functions for Kernel Smoothing Supporting Wand & Jones (1995)*. R package version 2.23-15, URL <https://CRAN.R-project.org/package=KernSmooth>.

- Wand MP, Ormerod JT (2008). “On Semiparametric Regression with O’Sullivan Penalized Splines.” *Australian and New Zealand Journal of Statistics*, **50**(2), 179–198. doi:10.1111/j.1467-842X.2008.00507.x.
- Wand MP, Ormerod JT, Padoan SA, Frühwirth R (2011). “Mean Field Variational Bayes for Elaborate Distributions.” *Bayesian Analysis*, **6**(4), 847–900. doi:10.1214/11-ba631.
- Wang SSJ, Wand MP (2011). “Using Infer.NET for Statistical Analyses.” *The American Statistician*, **65**(2), 115–126. doi:10.1198/tast.2011.10169.
- Winn J, Bishop CM (2005). “Variational Message Passing.” *Journal of Machine Learning Research*, **6**, 661–694.
- Wood SN (2003). “Thin Plate Regression Splines.” *Journal of the Royal Statistical Society B*, **65**(1), 95–114. doi:10.1111/1467-9868.00374.
- Wu L (2002). “A Joint Model for Nonlinear Mixed-Effects Models with Censored Response and Covariates Measured with Error, With Application to AIDS Studies.” *Journal of the American Statistical Association*, **97**(460), 700–709. doi:10.1198/016214502388618744.

A. Using Infer.NET to fit a simple semiparametric model

This section provides an extensive description of the Infer.NET program for semiparametric regression analysis of the onions data set based on model (16) in Section 3. The data are first transformed as explained in Section 3 and the transformed data are represented by y and $C = [X Z]$. Thereafter, the text files `K.txt`, `y.txt`, `Cmat.txt`, `sigsqBeta.txt`, `Au.txt`, `Aeps.txt` and `nIterVB.txt` are generated. The first three files contain the number of spline basis functions, y and C , respectively. The following three text files each contain a single number and represent the values for the hyperparameters: $\tau_\beta = 10^{-10}$, $A_\varepsilon = 10^5$ and $A_u = 10^5$. The last file, `nIterVB.txt`, contains a single positive number (i.e., 100) that specifies the number of mean field variational Bayes iterations. All these files were set up in R.

The following paragraphs explain the different commands in the C# script. Firstly, all required C# and Infer.NET libraries are loaded:

```
using System;
using System.Collections.Generic;
using System.Text;
using System.IO;
using System.Reflection;
using MicrosoftResearch.Infer;
using MicrosoftResearch.Infer.Distributions;
using MicrosoftResearch.Infer.Maths;
using MicrosoftResearch.Infer.Models;
```

The values for the hyperparameters, number of variational Bayes iterations and number of spline basis functions are imported using the commands:

```
double tauBeta = new DataTable(dir + "tauBeta.txt").DoubleNumeric;
double Au = new DataTable(dir + "Au.txt").DoubleNumeric;
double Aeps = new DataTable(dir + "Aeps.txt").DoubleNumeric;
int nIterVB = new DataTable(dir + "nIterVB.txt").IntNumeric;
int K = new DataTable(dir + "K.txt").IntNumeric;
```

Reading in the data y and C proceeds as follows:

```
DataTable yTable = new DataTable(dir + "y.txt");
DataTable C = new DataTable(dir + "Cmat.txt");
```

These commands make use of the C# object class ‘`DataTable`’. This class, which was written by the authors, facilitates the input of general rectangular arrays of numbers when available in a text file. The following line extracts the number of observations from the data matrix C :

```
int n = C.numRow;
```

The observed values of C and y , which are stored in objects `C` and `yTable`, are assigned to the objects `cvec` and `y` via the next set of commands:

```
Vector[] cvecTable = new Vector[n];
for (int i = 0; i < n; i++) {
```

```

    cvecTable[i] = Vector.Zero(3 + K);
    for (int j = 0; j < 3 + K; j++)
        cvecTable[i][j] = C.DataMatrix[i, j];
}
Range index = new Range(n).Named("index");
VariableArray<Vector> cvec = Variable.Array<Vector>(index).Named("cvec");
cvec.ObservedValue = cvecTable;
VariableArray<double> y = Variable.Array<double>(index).Named("y");
y.ObservedValue = yTable.arrayForIN;

```

The function `arrayForIN` is member of the class ‘`DataTable`’ and stores the data as an array to match the type of `y.ObservedValue`. The resulting objects are directly used to specify the prior distributions. First, the commands in code chunk (17) in Section 3 are used to set up priors for τ_u and τ_ε , while the trick involving the auxiliary variable $a \equiv 0$ in (16) is coded as:

```

Vector zeroVec = Vector.Zero(3 + K);
PositiveDefiniteMatrix betauPrecDummy =
    new PositiveDefiniteMatrix(3 + K, 3 + K);
for (int j = 0; j < 3 + K; j++)
    betauPrecDummy[j, j] = kappa;
Variable<Vector> betauWork =
    Variable.VectorGaussianFromMeanAndPrecision(zeroVec,
        betauPrecDummy).Named("betauWork");
Variable<double>[] a = new Variable<double>[3 + K];
for (int j = 0; j < 3; j++) {
    a[j] = Variable.GaussianFromMeanAndVariance(betauDummy[j], tauBeta);
    a[j].ObservedValue = 0.0;
}
for (int k = 0 ; k < K; k++) {
    a[3 + k] = Variable.GaussianFromMeanAndPrecision(betauDummy[3 + k], tauU);
    a[3 + k].ObservedValue = 0.0;
}

```

The command to specify the likelihood for model (16) is listed as code chunk (18) in Section 3. The inference engine and the number of mean field variational Bayes iterations are specified by means of code chunk (19). Finally, the following commands write the estimated values for the parameters of the approximate posteriors to files named `mu.q.betau.txt`, `Sigma.q.betau.txt`, `parms.q.tauEps.txt` and `parms.q.tauu.txt`:

```

SaveData.SaveTable(engine.Infer<VectorGaussian>(betauWork).GetMean(),
    dir + "mu.q.betau.txt");
SaveData.SaveTable(engine.Infer<VectorGaussian>(
    betauWork).GetVariance(), dir + "Sigma.q.betau.txt");
SaveData.SaveTable(engine.Infer<Gamma>(tauEps),
    dir + "parms.q.tauEps.txt");
SaveData.SaveTable(engine.Infer<Gamma>(tauu), dir + "parms.q.tauu.txt");

```

These commands involve the C# function named `SaveTable`, which is a method for the class ‘`SaveData`’. We wrote `SaveTable` to facilitate writing the output from *Infer.NET* to a text file.

Summary plots, such as Figure 2, can be made in R after back-transforming the approximate posterior density parameters.

Affiliation:

Jan Luts
TheSearchParty.com
Level 1, 79 Commonwealth Street
Surry Hills 2011, Australia
E-mail: jan.m.luts@gmail.com

Shen S. J. Wang
Department of Mathematics
The University of Queensland
Brisbane 4072, Australia
E-mail: gooney1980@gmail.com

John T. Ormerod
School of Mathematics and Statistics
University of Sydney
Sydney 2006, Australia
E-mail: jormerod@sydney.edu.au
URL: <http://www.maths.usyd.edu.au/u/jormerod/>

Matt P. Wand
School of Mathematical Sciences
University of Technology Sydney
Broadway 2007, Australia
E-mail: matt.wand@uts.edu.au
URL: <http://matt-wand.utsacademics.info/>