

Variational Message Passing for Elaborate Response Regression Models

M. W. McLean* and M. P. Wand†

Abstract. We build on recent work concerning message passing approaches to approximate fitting and inference for arbitrarily large regression models. The focus is on regression models where the response variable is modeled to have an elaborate distribution, which is loosely defined to mean a distribution that is more complicated than common distributions such as those in the Bernoulli, Poisson and Normal families. Examples of elaborate response families considered here are the Negative Binomial and t families. Variational message passing is more challenging due to some of the conjugate exponential families being non-standard and numerical integration being needed. Nevertheless, a factor graph fragment approach means the requisite calculations only need to be done once for a particular elaborate response distribution family. Computer code can be compartmentalized, including that involving numerical integration. A major finding of this work is that the modularity of variational message passing extends to elaborate response regression models.

Keywords: Bayesian computing, factor graph, generalized additive models, generalized linear mixed models, mean field variational Bayes, support vector machine classification.

MSC 2010 subject classifications: Primary 62F15, 62J05; secondary 62G08.

1 Introduction

We extend the variational message passing (VMP) body of work to accommodate elaborate response regression models. The notion of factor graph fragments, introduced in Wand (2017), is the vehicle for this extension. It affords a modular approach to mean field variational Bayes fitting and inference for large regression models. The factor graph fragment updates treated here only need to be derived and implemented once. Their addition to the variational message passing arsenal allows for fancier models, such as those having Negative Binomial and t responses, to be fitted.

VMP (Winn and Bishop, 2005; Minka, 2005; Minka and Winn, 2008) is a prescription for obtaining mean field variational Bayes approximations to posterior density functions that is amenable to modularization. The factor graph version of VMP (e.g Minka and Winn, 2008, Appendix A) is particularly attractive in this regard. Wand (2017) uses the notion of factor graph fragments to aid modularization for semiparametric regression

*School of Mathematical and Physical Sciences, University of Technology Sydney, P.O. Box 123, Broadway 2007, Australia, mathew.w.mclean@gmail.com

†School of Mathematical and Physical Sciences, University of Technology Sydney, P.O. Box 123, Broadway 2007, Australia, matt.wand@uts.edu.au

models – a large class of regression-type models that includes, for example, generalized linear mixed models, generalized additive models and varying coefficient models (e.g. Ruppert et al., 2003). However, the fragments in Wand (2017) only accommodate Gaussian, Bernoulli and Poisson response models. If, for example, a Negative Binomial response model is of interest then new fragment updates for this family are needed. Section 3.1 plugs this gap. Other elaborate response families are also treated in Section 3. Whilst we do not cover all possible families, our derivations for some elaborate families provide blueprints for future fragment derivations.

A major difference between simple response models and elaborate response models is that the latter involves non-standard exponential families. For the examples covered here four exponential families, beyond those covered in Wand (2017), emerge. Two of them seem to have little or no presence in the literature. The sufficient statistic expectations, which are needed for VMP updates, are not expressible in terms of common functions and require either evaluation of special functions, quadrature or continued fraction approximation.

The main contributions of this article may be summarized as follows:

1. If an analyst wants to build a mean field variational Bayes inference engine for arbitrarily large regression models then the message update formulae given in Section 3 allow for particular elaborate response families to be included;
2. The derivations in Section S.3 of the online supplement (McLean and Wand, 2018) show how such update formulae can be obtained for the examples given in Section 3. They also serve as a template for handling other elaborate response likelihoods not covered here.

All of our new methodology is within the realm of deterministic variational approximate inference, with intractable integrals evaluated via quadrature. An alternative route is to use Monte Carlo methods to approximate such integrals, known as *stochastic* variational inference (e.g. Hoffman et al., 2013; Kucukelbir et al., 2017). See, for example, Titsias and Lázaro-Gredilla (2014) on the use of stochastic variational inference for non-conjugate circumstances similar to those arising in this article.

Some background on VMP is given in Section 2. Section 3 is the article’s centerpiece and gives the fragment update for six elaborate response likelihoods. Illustration of their utility is then provided in Section 4. Closing remarks are given in Section 5. Derivational details are given in an online supplement.

2 Variational Message Passing and Factor Graph Fragments

Variational message passing (VMP) is an approach to obtaining mean field variational Bayes approximate posterior density functions in potentially large graphical models. It uses the concept of *message passing* on a *factor graph*.

Our starting point is a Bayesian statistical model with observed data \mathbf{D} and parameter vector $\boldsymbol{\theta}$. The posterior density function $p(\boldsymbol{\theta}|\mathbf{D})$ is usually analytically intractable and a *mean field variational* approximation $q^*(\boldsymbol{\theta})$ to $p(\boldsymbol{\theta}|\mathbf{D})$ is the minimizer of the Kullback–Leibler divergence

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

subject to the product density restriction $q(\boldsymbol{\theta}) = \prod_{i=1}^M q(\boldsymbol{\theta}_i)$ where $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$ is some partition of $\boldsymbol{\theta}$. The optimal q -density functions can be shown to satisfy

$$q^*(\boldsymbol{\theta}_i) \propto E_{q(\boldsymbol{\theta} \setminus \boldsymbol{\theta}_i)} \{p(\boldsymbol{\theta}_i|\mathbf{D}, \boldsymbol{\theta} \setminus \boldsymbol{\theta}_i)\}, \quad 1 \leq i \leq M, \tag{1}$$

where $\boldsymbol{\theta} \setminus \boldsymbol{\theta}_i$ denotes the entries of $\boldsymbol{\theta}$ with $\boldsymbol{\theta}_i$ omitted. Expression (1) gives rise to an iterative scheme for determination of the optimal parameters of the $q^*(\boldsymbol{\theta}_i)$, which is known as *mean field variational Bayes*. A listing of such a scheme is provided by Algorithm 1 of Ormerod and Wand (2010).

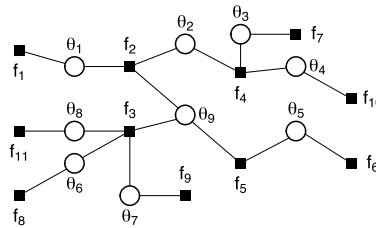


Figure 1: Factor graph representation of the dependence of the stochastic nodes $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_9$ on the factors f_1, \dots, f_{11} for the example given by (2).

VMP arrives at the same approximation via message passing on an appropriate factor graph. Figure 1 is an example factor graph corresponding to an $M = 9$ example with

$$p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_9, \mathbf{D}) = f_1(\boldsymbol{\theta}_1) f_2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_9) f_3(\boldsymbol{\theta}_6, \boldsymbol{\theta}_7, \boldsymbol{\theta}_8, \boldsymbol{\theta}_9) f_4(\boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \boldsymbol{\theta}_4) f_5(\boldsymbol{\theta}_5, \boldsymbol{\theta}_9) \tag{2} \\ \times f_6(\boldsymbol{\theta}_5) f_7(\boldsymbol{\theta}_3) f_8(\boldsymbol{\theta}_6) f_9(\boldsymbol{\theta}_7) f_{10}(\boldsymbol{\theta}_4) f_{11}(\boldsymbol{\theta}_8).$$

At least one of the f_j involves the data vector \mathbf{D} , but this dependence is suppressed. The unshaded circles are called *stochastic nodes* and the shaded rectangles are the *factors*. The word *node* is used for either a stochastic node or a factor and two nodes are neighbors of each other if they are joined by an edge. The edges join factors to stochastic nodes that are included in that factor. The θ_i indices connected to the j th factor are denoted by $\text{neighbors}(j)$. For example, $\text{neighbors}(3) = \{6, 7, 8, 9\}$. Fuller details are in Sections 2.4 and 2.5 of Wand (2017).

A *message* passed between any two neighboring nodes is a particular function of the stochastic node that either sends or receives the message. Rather than using (1), the optimal q -densities are obtained from

$$q^*(\boldsymbol{\theta}_i) \propto \prod_{j:i \in \text{neighbors}(j)} m_{f_j \rightarrow \boldsymbol{\theta}_i}^*(\boldsymbol{\theta}_i), \quad (3)$$

where the $m_{f_j \rightarrow \boldsymbol{\theta}_i}^*(\boldsymbol{\theta}_i)$ are the optimal messages passed to $\boldsymbol{\theta}_i$ from each of the factors f_j in $p(\boldsymbol{\theta}, \mathbf{D})$ that involve $\boldsymbol{\theta}_i$. For each j , this subset of $\{1, \dots, M\}$ is denoted by $\text{neighbors}(j)$ due to the definition of a factor graph, in which an edge is drawn between the $\boldsymbol{\theta}_i$ and f_j nodes if and only if f_j depends on $\boldsymbol{\theta}_i$.

Letting N denote the number of factors, for each $1 \leq i \leq M$ and $1 \leq j \leq N$ the VMP stochastic node to factor message updates are

$$m_{\boldsymbol{\theta}_i \rightarrow f_j}(\boldsymbol{\theta}_i) \leftarrow \propto \prod_{j' \neq j: i \in \text{neighbors}(j')} m_{f_{j'} \rightarrow \boldsymbol{\theta}_i}(\boldsymbol{\theta}_i) \quad (4)$$

and the factor to stochastic node message updates are

$$m_{f_j \rightarrow \boldsymbol{\theta}_i}(\boldsymbol{\theta}_i) \leftarrow \propto \exp \left[E_{f_j \rightarrow \boldsymbol{\theta}_i} \left\{ \log f_j(\boldsymbol{\theta}_{\text{neighbors}(j)}) \right\} \right], \quad (5)$$

where $E_{f_j \rightarrow \boldsymbol{\theta}_i}$ denotes expectation with respect to the density function

$$\frac{\prod_{i' \in \text{neighbors}(j) \setminus \{i\}} m_{f_j \rightarrow \boldsymbol{\theta}_{i'}}(\boldsymbol{\theta}_{i'}) m_{\boldsymbol{\theta}_{i'} \rightarrow f_j}(\boldsymbol{\theta}_{i'})}{\prod_{i' \in \text{neighbors}(j) \setminus \{i\}} \int m_{f_j \rightarrow \boldsymbol{\theta}_{i'}}(\boldsymbol{\theta}_{i'}) m_{\boldsymbol{\theta}_{i'} \rightarrow f_j}(\boldsymbol{\theta}_{i'}) d\boldsymbol{\theta}_{i'}}. \quad (6)$$

In (4) and (5) the $\leftarrow \propto$ symbol means that the function of $\boldsymbol{\theta}_i$ on the left-hand side is updated according to the expression on the right-hand side but that multiplicative factors not depending on $\boldsymbol{\theta}_i$ can be ignored. If $\text{neighbors}(j) \setminus \{i\} = \emptyset$ then the expectation in (5) can be dropped and the right-hand side of (5) is proportional to $f_j(\boldsymbol{\theta}_{\text{neighbors}(j)})$.

VMP fitting involves iteration of the updates (4) and (5)–(6) over each of the factors until the changes in all messages are negligible. When convergence is reached, the optimal q -densities of the model parameters are obtained from (3).

The algebra and coding for VMP can be compartmentalized using the notion of *factor graph fragments*, or *fragments* for short.

Definition. A *factor graph fragment*, or *fragment* for short, is a sub-graph of a factor graph consisting of a single factor and each of the stochastic nodes that are neighbors of the factor.

In the context of the current article, the fragment approach means that switching from a large regression-type model with a Gaussian likelihood to one with, say, a t likelihood can be achieved by replacing the Gaussian likelihood fragment by t likelihood fragments. The remainder of the model is unaffected in terms of the VMP updates.

Table 1 of Wand (2017) lists five fragments that are fundamental to semiparametric regression analysis via VMP. As explained there, a wide range of semiparametric regression models are accommodated by these five fragments but only for the Gaussian

response case. In Section 5 of Wand (2017), additional fragments are introduced to handle logistic, probit and Poisson regression models. The next section adds to these response fragments.

3 Fragment Updates for Elaborate Response Likelihoods

We now provide fragment updates that allow for six more response distributions to be handled within the VMP framework. Most of them may be viewed as elaborations of the likelihoods covered by Wand (2017). For example, the Negative Binomial likelihood extends the Poisson likelihood for count response data and the t and Skew Normal likelihoods extend the Gaussian likelihood in different ways.

Each of the elaborate response likelihoods considered in this section are re-expressed in terms of auxiliary variables and more common distributions. This affords tractability, but comes at the cost of less accuracy compared with the case where auxiliary variables are not introduced. The auxiliary variables route is driven by the practical advantages of message updates being either closed form or requiring only univariate numerical integration. The alternative route, without auxiliary variables, is much more numerically challenging and often impractical.

Table 1 provides details on each of the distributions used in this article. It uses the following notation for the $N(0, 1)$ density and cumulative distribution functions:

$$\phi(x) \equiv (2\pi)^{-1/2} \exp(-\frac{1}{2} x^2) \quad \text{and} \quad \Phi(x) \equiv \int_{-\infty}^x \phi(t) dt.$$

An additional functional notation is digamma(x) $\equiv \frac{d}{dx} \log\{\Gamma(x)\}$.

For a vector \mathbf{a} and scalar function s we let $s(\mathbf{a})$ denote the vector containing the element-wise evaluations of s . Also, $\mathbf{A} \odot \mathbf{B}$ and \mathbf{A}/\mathbf{B} respectively denote the element-wise product and element-wise quotient of vectors \mathbf{A} and \mathbf{B} having the same sizes. If \mathbf{A} is a $d \times d$ matrix then $\text{vec}(\mathbf{A})$ is the $d^2 \times 1$ vector obtained by stacking the columns of \mathbf{A} underneath each other in order from left to right. If \mathbf{a} is a $d^2 \times 1$ vector then $\text{vec}^{-1}(\mathbf{a})$ is the $d \times d$ matrix formed from listing the entries of \mathbf{a} in column-wise fashion in order from left to right. The $d \times 1$ vector containing the diagonal entries of a $d \times d$ matrix \mathbf{A} is denoted by $\text{diagonal}(\mathbf{A})$.

The $d \times 1$ vector $\mathbf{1}_d$ is such that all of its entries are equal to 1. The $d \times 1$ vector \mathbf{e}_i is such that its i th entry is equal to 1 and all other entries are zero.

For a $d \times 1$ vector \mathbf{v}_1 and a $d^2 \times 1$ vector \mathbf{v}_2 such that $\text{vec}^{-1}(\mathbf{v}_2)$ is symmetric we define:

$$G_{\text{VMP}} \left(\begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}; \mathbf{Q}, \mathbf{r}, s \right) \equiv -\frac{1}{8} \text{tr} \left(\mathbf{Q} \{ \text{vec}^{-1}(\mathbf{v}_2) \}^{-1} [\mathbf{v}_1 \mathbf{v}_1^T \{ \text{vec}^{-1}(\mathbf{v}_2) \}^{-1} - 2\mathbf{I}] \right) - \frac{1}{2} \mathbf{r}^T \{ \text{vec}^{-1}(\mathbf{v}_2) \}^{-1} \mathbf{v}_1 - \frac{1}{2} s.$$

The secondary arguments of G_{VMP} are a $d \times d$ matrix \mathbf{Q} , a $d \times 1$ vector \mathbf{r} and $s \in \mathbb{R}$. The genesis of the G_{VMP} function is the fact that

distribution	density/probability function in x	abbreviation
Multinomial	$\prod_{k=1}^K \pi_k^{x_k}; x_k = 0, 1, 1 \leq k \leq K; \sum_{k=1}^K \pi_k = 1$	Multinomial($1, \boldsymbol{\pi}$)
Poisson	$\lambda^x e^{-\lambda}/x!; x = 0, 1, \dots; \lambda > 0$	Poisson(λ)
Negative Binomial	$\frac{\kappa^\kappa \Gamma(x + \kappa) \mu^x}{\Gamma(\kappa)(\kappa + \mu) \Gamma(x + 1)}; x = 0, 1, \dots; \kappa, \mu > 0$	Negative-Binomial(μ, κ)
t	$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu} \Gamma(\nu/2) [1 + \{(x - \mu)/\sigma\}^2/\nu]^{\frac{\nu+1}{2}}}; \sigma, \nu > 0$	$t(\mu, \sigma, \nu)$
Asymmetric Laplace	$\frac{\tau(1-\tau)}{\sigma} \exp[-\frac{1}{2} \frac{x-\mu}{\sigma} + (\tau - \frac{1}{2}) (\frac{x-\mu}{\sigma})]; \sigma > 0, 0 < \tau < 1$	Asymmetric-Laplace(μ, σ, τ)
Skew Normal	$\frac{2}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) \Phi\left(\frac{\lambda(x - \mu)}{\sigma}\right); \sigma > 0$	Skew-Normal(μ, σ, λ)
Finite Normal Mixture	$\sum_{k=1}^K \frac{w_k}{\sigma s_k} \phi\left(\frac{(x - \mu)/\sigma - m_k}{s_k}\right); w_k, s_k > 0, \sum_{k=1}^K w_k = 1$	Normal-Mixture($\mu, \sigma, \mathbf{w}, \mathbf{m}, \mathbf{s}$)
Gamma	$\frac{B^A x^{A-1} e^{-Bx}}{\Gamma(A)}; x > 0, A, B > 0$	Gamma(A, B)
Inverse- χ^2	$\frac{(\lambda/2)^{\kappa/2} x^{-(\kappa/2)-1} e^{-(\lambda/2)/x}}{\Gamma(\kappa/2)}; x > 0; \kappa, \lambda > 0$	Inverse- $\chi^2(\kappa, \lambda)$

Table 1: Distributions used in this article and their corresponding density/probability functions.

$$E_{\boldsymbol{\theta}}\{-\frac{1}{2}(\boldsymbol{\theta}^T \mathbf{Q} \boldsymbol{\theta} - 2\mathbf{r}^T \boldsymbol{\theta} + s)\} = G_{\text{VMP}}(\boldsymbol{\eta}; \mathbf{Q}, \mathbf{r}, s),$$

when $\boldsymbol{\theta}$ is a $d \times 1$ Multivariate Normal random vector with natural parameter vector $\boldsymbol{\eta}$. A last piece of notation is

$$\boldsymbol{\eta}_{f \leftrightarrow \boldsymbol{\theta}} \equiv \boldsymbol{\eta}_{f \rightarrow \boldsymbol{\theta}} + \boldsymbol{\eta}_{\boldsymbol{\theta} \rightarrow f}$$

for any natural parameter $\boldsymbol{\eta}$, factor f and stochastic node $\boldsymbol{\theta}$.

3.1 Negative Binomial Likelihood

The Negative Binomial likelihood fragments are concerned with the likelihood specification

$$y_i | \boldsymbol{\theta}, \kappa \stackrel{\text{ind.}}{\sim} \text{Negative-Binomial}[\exp\{(\mathbf{A}\boldsymbol{\theta})_i\}, \kappa], \quad 1 \leq i \leq n. \quad (7)$$

Introduce Gamma auxiliary random variables $a_i | \boldsymbol{\theta}, \kappa \stackrel{\text{ind.}}{\sim} \text{Gamma}[\kappa, \kappa \exp\{-(\mathbf{A}\boldsymbol{\theta})_i\}]$, $1 \leq i \leq n$. Then standard distribution theoretical manipulations lead to (7) being equivalent to

$$y_i | a_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(a_i), \quad a_i | \boldsymbol{\theta}, \kappa \stackrel{\text{ind.}}{\sim} \text{Gamma}[\kappa, \kappa \exp\{-(\mathbf{A}\boldsymbol{\theta})_i\}].$$

The relevant factor graph fragments are shown in Figure 2 and corresponds to the mean field restriction

$$q(\boldsymbol{\theta}, \kappa, \mathbf{a}) = q(\boldsymbol{\theta})q(\kappa) \left\{ \prod_{i=1}^n q(a_i) \right\}$$

that was used in Luts and Wand (2015).

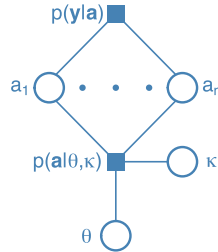


Figure 2: Fragments for the Negative Binomial likelihood specification with independent Gamma auxiliary variables a_1, \dots, a_n .

First note that

$$m_{p(\mathbf{a}|\boldsymbol{\theta}, \kappa) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) = \exp \left[- E_{q(\kappa)}(\kappa) \{ \mathbf{1}_n^T \mathbf{A} \boldsymbol{\theta} + E_{q(\mathbf{a})}(\mathbf{a})^T \exp(-\mathbf{A} \boldsymbol{\theta}) \} \right], \quad (8)$$

which is not conjugate with Multivariate Normal messages passed to $\boldsymbol{\theta}$ from other factors. Instead, we replace (8) with

$$\tilde{m}_{p(\mathbf{a}|\boldsymbol{\theta}, \kappa) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) \equiv \exp \left\{ \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta}\boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{a}|\boldsymbol{\theta}, \kappa) \rightarrow \boldsymbol{\theta}} \right\} \quad (9)$$

to enforce conjugacy with Multivariate Normal messages. This is an instance of non-conjugate VMP (Knowles and Minka, 2011). We assume that each of the messages that $\boldsymbol{\theta}$ receives from factors outside of the Negative Binomial likelihood fragments are within the Multivariate Normal family. This leads to $q^*(\boldsymbol{\theta})$ having a Multivariate Normal distribution.

As explained in Section S.3.1 of the online supplement, the message from $p(\mathbf{a}|\boldsymbol{\theta}, \kappa)$ to κ takes the form

$$m_{p(\mathbf{a}|\boldsymbol{\theta}, \kappa) \rightarrow \kappa}(\kappa) = \exp \left\{ \left[\begin{array}{c} \kappa \log(\kappa) - \log\{\Gamma(\kappa)\} \\ \kappa \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{a}|\boldsymbol{\theta}, \kappa) \rightarrow \kappa} \right\},$$

which is proportional to the *Moon Rock* exponential family of density functions described in Section S.2.4 of the online supplement. We assume messages passed to κ from factors outside of the Negative Binomial likelihood fragments are also within the

Algorithm 1 The inputs, updates and outputs of the Negative Binomial likelihood fragment.

Data Inputs: y, A .

Parameter Inputs: $\boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow p(a|\boldsymbol{\theta}, \kappa), \boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \kappa, \boldsymbol{\eta}_{\kappa} \rightarrow p(a|\boldsymbol{\theta}, \kappa)$.

Updates:

$$\mu_{q(\kappa)} \leftarrow (ET)_2^{\text{MR}}(\boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \leftrightarrow \kappa)$$

$$\boldsymbol{\omega}_1 \leftarrow -\frac{1}{2} \mathbf{A} \left\{ \text{vec}^{-1} \left((\boldsymbol{\eta}_{p(a|\boldsymbol{\theta})} \leftrightarrow \boldsymbol{\theta})_2 \right) \right\}^{-1} (\boldsymbol{\eta}_{p(a|\boldsymbol{\theta})} \leftrightarrow \boldsymbol{\theta})_1$$

$$\boldsymbol{\omega}_2 \leftarrow \exp \left(-\boldsymbol{\omega}_1 - \frac{1}{4} \text{diagonal} \left[\mathbf{A} \left\{ \text{vec}^{-1} \left((\boldsymbol{\eta}_{p(a|\boldsymbol{\theta})} \leftrightarrow \boldsymbol{\theta})_2 \right) \right\}^{-1} \mathbf{A}^T \right] \right)$$

$$\boldsymbol{\omega}_3 \leftarrow \{ \boldsymbol{\omega}_2 \odot (\mathbf{y} + \mu_{q(\kappa)} \mathbf{1}_n) \} / (\mathbf{1}_n + \mu_{q(\kappa)} \boldsymbol{\omega}_2)$$

$$\boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \boldsymbol{\theta} \leftarrow \mu_{q(\kappa)} \begin{bmatrix} \mathbf{A}^T \{ \boldsymbol{\omega}_3 \odot (\mathbf{1}_n + \boldsymbol{\omega}_1) - \mathbf{1}_n \} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_3) \mathbf{A}) \end{bmatrix}$$

$$\boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \kappa \leftarrow \begin{bmatrix} n \\ \mathbf{1}_n^T \{ \text{digamma}(\mu_{q(\kappa)} \mathbf{1}_n + \mathbf{y}) - \boldsymbol{\omega}_1 \\ -\log(\mathbf{1}_n + \mu_{q(\kappa)} \boldsymbol{\omega}_2) - \boldsymbol{\omega}_3 \} \end{bmatrix}$$

Parameter Outputs: $\boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \kappa$.

Moon Rock family or at least conjugate with the Moon Rock family. For example, if the only other factor passing messages to κ is its prior density function $p(\kappa)$ then we require that $p(\kappa)$ is a Moon Rock density function or conjugate with one. Note that, for example, Exponential density functions (Gamma(1, B) density functions in the notation of Table 1) are conjugate with respect to the Moon Rock family but, strictly speaking, not within the Moon Rock family since $\alpha = 0$ in the notation of Section S.2.4. Hence, setting

$$p(\kappa) = B \exp(-B \kappa), \quad \kappa > 0,$$

for any $B > 0$ is permissible under the conjugacy constraint since it implies that

$$m_{p(\kappa) \rightarrow \kappa}(\kappa) = \exp \left\{ \left[\begin{array}{c} \kappa \log(\kappa) - \log\{\Gamma(\kappa)\} \\ \kappa \end{array} \right]^T \left[\begin{array}{c} 0 \\ -B \end{array} \right] \right\}.$$

which is conjugate with respect to $m_{p(a|\boldsymbol{\theta}, \kappa)} \rightarrow \kappa(\kappa)$.

Algorithm 1 lists the inputs, updates and outputs for the Negative Binomial likelihood fragments. The derivations are given in Section S.3.1 of the online supplement.

The $(ET)_2^{\text{MR}}$ notation, used in the first update, is explained in Section S.2.4 of the online supplement.

In Section 4.2 we provide illustration of Algorithm 1 in the context of additive model analysis.

3.2 t Likelihood

The t -distribution likelihood fragments arise from the likelihood specification

$$y_i | \boldsymbol{\theta}, \sigma, \nu \stackrel{\text{ind.}}{\sim} t\left((\mathbf{A}\boldsymbol{\theta})_i, \sigma, \nu\right), \quad 1 \leq i \leq n. \tag{10}$$

This likelihood is frequently used in regression applications as a robustness mechanism (e.g. Lange et al., 1989). If we introduce Inverse- χ^2 auxiliary random variables $a_i \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(\nu, \nu)$, $1 \leq i \leq n$, then (10) is equivalent to

$$y_i | \boldsymbol{\theta}, \sigma^2, a_i \stackrel{\text{ind.}}{\sim} N\left((\mathbf{A}\boldsymbol{\theta})_i, a_i \sigma^2\right), \quad a_i | \nu \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(\nu, \nu). \tag{11}$$

It is common to use this representation of the t distribution for Bayesian computing. For example, the Markov chain Monte Carlo scheme of Verdinelli and Wasserman (1991) and the mean field variational Bayes scheme of Tipping and Lawrence (2003) each rely upon (11).

Figure 3 shows the factor graph fragments for the auxiliary variable representation (11) with q -density product restriction

$$q(\boldsymbol{\theta}, \sigma^2, \nu, \mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2)q(\nu) \left\{ \prod_{i=1}^n q(a_i) \right\}.$$

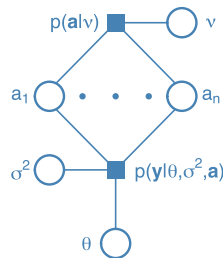


Figure 3: Fragments for the t likelihood specification with the shape parameter prior with independent Inverse- $\chi^2(\nu, \nu)$ auxiliary variables a_1, \dots, a_n .

The message from $p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2, \mathbf{a})$ to $\boldsymbol{\theta}$ is proportional to a Multivariate Normal density function, while that from $p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2, \mathbf{a})$ to σ^2 is within the Inverse- χ^2 family.

The message from $p(\mathbf{a}|\nu)$ to ν has the form

$$m_{p(\mathbf{a}|\nu) \rightarrow \nu}(\nu) = \exp \left\{ \left[\begin{array}{c} (\nu/2) \log(\nu/2) - \log\{\Gamma(\nu/2)\} \\ \nu/2 \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{a}|\nu) \rightarrow \nu} \right\}$$

with details given in Section S.3.2 of the online supplement. Note that $m_{p(\mathbf{a}|\nu) \rightarrow \nu}(\nu)$ is proportional to a factor of 2 rescaling of the Moon Rock exponential family of density functions introduced in Section S.2.4 of the online supplement. The conjugacy constraint dictates that

$$m_{\nu \rightarrow p(\mathbf{a}|\nu)}(\nu) = \exp \left\{ \left[\begin{array}{c} (\nu/2) \log(\nu/2) - \log\{\Gamma(\nu/2)\} \\ \nu/2 \end{array} \right]^T \boldsymbol{\eta}_{\nu \rightarrow p(\mathbf{a}|\nu)} \right\},$$

which occurs if all message passed to ν from factors outside of the t likelihood fragments are also within the same rescaled Moon Rock family, or at least conjugate with respect to it. The $(\mathbf{E}\mathbf{T})_2^{\text{MR}}$ notation is defined in Section S.2.4 of the online supplement.

Algorithm 2 provides the inputs, updates and outputs for the t likelihood fragments. The derivations are given in Section S.3.2 of the online supplement.

3.3 Asymmetric Laplace Likelihood

Now consider the Asymmetric Laplace likelihood specification

$$y_i | \boldsymbol{\theta}, \sigma^2 \stackrel{\text{ind.}}{\sim} \text{Asymmetric-Laplace}((\mathbf{A}\boldsymbol{\theta})_i, \sigma, \tau), \quad 1 \leq i \leq n, \quad (12)$$

where $0 < \tau < 1$ is a fixed constant. As explained in, for example, Yu and Moyeed (2001), the likelihood specification (12) corresponds to τ th-quantile regression. Yang et al. (2016) discuss valid posterior inference for Bayesian quantile regression.

If we introduce auxiliary random variables $a_i \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(2, 1)$, $1 \leq i \leq n$, then Proposition 3.2.1 of Kotz et al. (2001) implies that (12) is equivalent to

$$y_i | \boldsymbol{\theta}, \sigma^2, \mathbf{a} \stackrel{\text{ind.}}{\sim} N \left((\mathbf{A}\boldsymbol{\theta})_i + \frac{(\frac{1}{2} - \tau)\sigma}{a_i\tau(1-\tau)}, \frac{\sigma^2}{a_i\tau(1-\tau)} \right), \quad a_i \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(2, 1). \quad (13)$$

We assume that the optimal q -density admits the product restriction

$$q(\boldsymbol{\theta}, \sigma^2, \mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2)q(\mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2) \prod_{i=1}^n q(a_i).$$

The corresponding factor graph fragments are shown in Figure 4.

As shown in Section S.3.3 of the online supplement,

$$m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) = \exp \left\{ \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta}\boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \rightarrow \boldsymbol{\theta}} \right\},$$

Algorithm 2 The inputs, updates and outputs of the t likelihood fragment.

Data Inputs: y, \mathbf{A} .

Parameter Inputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}), \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2,$

$$\boldsymbol{\eta}_{\sigma^2} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}), \boldsymbol{\eta}_{p(\mathbf{a}|\nu)} \rightarrow \nu, \boldsymbol{\eta}_{\nu} \rightarrow p(\mathbf{a}|\nu).$$

Updates:

$$\mu_{q(1/\sigma^2)} \leftarrow \left\{ (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2)_1 + 1 \right\} / (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2)_2$$

$$\mu_{q(\nu)} \leftarrow 2(ET)_2^{\text{MR}}(\boldsymbol{\eta}_{p(\mathbf{a}|\nu)} \leftrightarrow \nu)$$

$$\boldsymbol{\omega}_4 \leftarrow \left[G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{A}, \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{y}, y_i^2 \right) \right]_{1 \leq i \leq n}$$

$$\boldsymbol{\omega}_5 \leftarrow \frac{(\mu_{q(\nu)} + 1) \mathbf{1}_n}{\mu_{q(\nu)} \mathbf{1}_n - 2\mu_{q(1/\sigma^2)} \boldsymbol{\omega}_4}$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta} \leftarrow \mu_{q(1/\sigma^2)} \left[\begin{array}{c} \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_5) \mathbf{y} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_5) \mathbf{A}) \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2 \leftarrow \left[\begin{array}{c} -\frac{1}{2} n \\ G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_5) \mathbf{A}, \right. \\ \left. \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_5) \mathbf{y}, \mathbf{y}^T \text{diag}(\boldsymbol{\omega}_5) \mathbf{y} \right) \end{array} \right]$$

$$\boldsymbol{\eta}_{p(\mathbf{a}|\nu)} \rightarrow \nu \leftarrow \left[\begin{array}{c} n \\ n \text{digamma} \left(\frac{\mu_{q(\nu)} + 1}{2} \right) - \mathbf{1}_n^T \left\{ \log \left(\frac{1}{2} \mu_{q(\nu)} \mathbf{1}_n - \mu_{q(1/\sigma^2)} \boldsymbol{\omega}_4 \right) \right. \\ \left. + \boldsymbol{\omega}_5 \right\} \end{array} \right]$$

Parameter Outputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2, \boldsymbol{\eta}_{p(\mathbf{a}|\nu)} \rightarrow \nu.$

which is conjugate with Multivariate Normal messages passed to $\boldsymbol{\theta}$ from factors outside of the Asymmetric Laplace likelihood fragments.

However, the message from the likelihood factor to σ^2 takes the form

$$m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2(\sigma^2) = \exp \left\{ \left[\begin{array}{c} \log(\sigma^2) \\ 1/\sigma \\ 1/\sigma^2 \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2 \right\},$$

which is not within a standard exponential family. However, $m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2(\sigma^2)$ is proportional to the family of density functions of random variables such that their reciprocal square roots are distributed according to members of a family proposed in

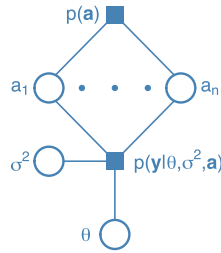


Figure 4: Fragments for the Asymmetric Laplace likelihood specification with independent Inverse- $\chi^2(2, 1)$ auxiliary variables a_1, \dots, a_n .

Nadarajah (2008). Sections S.2.2 and S.2.3 of the online supplement contain the relevant details. We will assume that messages passed to σ^2 from factors outside of the Asymmetric Laplace likelihood fragments are within the *Inverse Square Root Nadarajah* family (Section S.2.3 of the online supplement). Note that messages proportional to Inverse Chi-Squared density functions are conjugate with this family.

Algorithm 3 provides the inputs, updates and outputs for the Asymmetric Laplace likelihood fragments with derivations deferred to Section S.3.3 of the online supplement. Note that the second update of Algorithm 3 involves the function \mathcal{R}_ν , which is defined in Section S.1.2 of the online supplement. Efficient and stable computation of \mathcal{R}_ν is discussed there.

In Section 4.1 we show that Algorithm 3 facilitates quantile nonparametric regression embellishment of ordinary nonparametric regression.

Laplace Likelihood Special Case

The case of $\tau = \frac{1}{2}$ corresponds to the special case of the Laplace likelihood, and (12) reduces to *median* regression. In this special case, the second entry of $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \rightarrow \sigma^2}$ is zero and messages passed to σ^2 are proportional to Inverse Chi-Squared density functions. In addition, the $\mu_{q(1/\sigma)}$ update in Algorithm 3 is not needed and that for $\mu_{q(1/\sigma^2)}$ reduces to

$$\mu_{q(1/\sigma^2)} \leftarrow \left\{ (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \leftrightarrow \sigma^2})_1 + 1 \right\} / (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \leftrightarrow \sigma^2})_2,$$

where $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}) \leftrightarrow \sigma^2}$ is an Inverse Chi-Squared natural parameter vector.

3.4 Skew Normal Likelihood

In this section, we consider fragments involving the Skew Normal likelihood:

$$y_i | \boldsymbol{\theta}, \sigma^2, \lambda \sim \text{Skew-Normal}\{(\mathbf{A}\boldsymbol{\theta})_i, \sigma, \lambda\}, \quad 1 \leq i \leq n. \quad (14)$$

Regression-type models having Skew Normal responses may be found in, for example, Frühwirth-Schnatter and Pyne (2010) and Lachos et al. (2010).

Algorithm 3 The inputs, updates and outputs of the Asymmetric Laplace likelihood fragments.

Data Inputs: y, \mathbf{A}, τ .

Parameter Inputs: $\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a}), \boldsymbol{\eta}_{\sigma^2} \rightarrow p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a}),$
 $\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2.$

Updates:

$$\mu_{q(1/\sigma)} \leftarrow (ET)_2^{\text{ISRN}} \left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2 \right)$$

$$\mu_{q(1/\sigma^2)} \leftarrow (ET)_3^{\text{ISRN}} \left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2 \right)$$

$$\boldsymbol{\omega}_7 \leftarrow \left[G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{A}, \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{y}, y_i^2 \right) \right]_{1 \leq i \leq n}$$

$$\boldsymbol{\omega}_8 \leftarrow \left\{ -8 \tau^2 (1 - \tau)^2 \mu_{q(1/\sigma^2)} \boldsymbol{\omega}_7 \right\}^{-1/2}$$

$$\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta} \leftarrow \tau(1 - \tau) \mu_{q(1/\sigma^2)} \begin{bmatrix} \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_8) \mathbf{y} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_8) \mathbf{A}) \end{bmatrix}$$

$$+ \left(\tau - \frac{1}{2} \right) \mu_{q(1/\sigma)} \begin{bmatrix} \mathbf{A}^T \mathbf{1}_n \\ \mathbf{0} \end{bmatrix}$$

$$\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2 \leftarrow \begin{bmatrix} -n/2 \\ \left(\frac{1}{2} - \tau \right) \left[\mathbf{y} + \frac{1}{2} \mathbf{A} \left\{ \text{vec}^{-1} \left(\left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta} \right)_2 \right) \right\}^{-1} \right. \\ \left. \times \left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta} \right)_1 \right]^T \mathbf{1}_n \\ \left. \tau(1 - \tau) G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_8) \mathbf{A}, \right. \right. \\ \left. \left. \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_8) \mathbf{y}, \mathbf{y}^T \text{diag}(\boldsymbol{\omega}_8) \mathbf{y} \right) \right] \end{bmatrix}$$

Parameter Outputs: $\boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{p(y|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2.$

If we introduce auxiliary random variables $a_i \stackrel{\text{ind.}}{\sim} N(0, 1)$, $1 \leq i \leq n$, then Proposition 3 of Azzalini and Dalla Valle (1996) implies that (14) is equivalent to

$$y_i | \boldsymbol{\theta}, \sigma^2, \lambda, a_i \stackrel{\text{ind.}}{\sim} N \left((\mathbf{A}\boldsymbol{\theta})_i + \frac{\sigma \lambda |a_i|}{\sqrt{1 + \lambda^2}}, \frac{\sigma^2}{1 + \lambda^2} \right), \quad a_i \stackrel{\text{ind.}}{\sim} N(0, 1). \quad (15)$$

We assume the optimal q -density admits the product restriction

$$q(\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2)q(\lambda)q(\mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2)q(\lambda) \prod_{i=1}^n q(a_i).$$

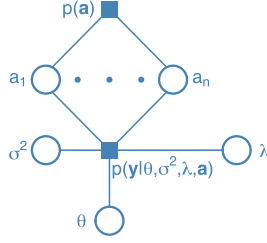


Figure 5: Fragments for the Skew Normal likelihood specification with independent $N(0, 1)$ auxiliary variables a_1, \dots, a_n .

The corresponding factor graph fragments are shown in Figure 5.

The messages passed from the likelihood factor to θ and σ^2 take the forms

$$m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \boldsymbol{\theta}}(\boldsymbol{\theta}) = \exp \left\{ \left[\begin{array}{c} \boldsymbol{\theta} \\ \text{vec}(\boldsymbol{\theta}\boldsymbol{\theta}^T) \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \boldsymbol{\theta}} \right\}$$

and

$$m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \sigma^2}(\sigma^2) = \exp \left\{ \left[\begin{array}{c} \log(\sigma^2) \\ 1/\sigma \\ 1/\sigma^2 \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \sigma^2} \right\}.$$

As for the Asymmetric Laplace likelihood fragments, the latter is within the Inverse Square Root Nadarajah family. The imposition of conjugacy means that we assume that all messages passed to σ^2 from factors outside of the Skew Normal likelihood fragments are also proportional to Inverse Square Root Nadarajah density functions.

The message from the likelihood factor to λ has the exponential family form

$$m_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \lambda}(\lambda) = \exp \left\{ \left[\begin{array}{c} \log(1 + \lambda^2) \\ \lambda^2 \\ \lambda\sqrt{1 + \lambda^2} \end{array} \right]^T \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}) \rightarrow \lambda} \right\}.$$

We have not been able to find any mention of this family in the literature. In Section S.2.5 of the online supplement we dub it the *Sea Sponge* family. We assume that each of the messages that λ receives from factors outside of the Skew Normal likelihood fragments are also proportional to Sea Sponge density functions. As an example, suppose that the only other factor that sends a message to λ is the prior density function $p(\lambda)$. Then, $m_{p(\lambda) \rightarrow \lambda}(\lambda) = p(\lambda)$ and, under conjugacy, $p(\lambda)$ must be of the form

$$p(\lambda) \propto \exp \left\{ \left[\begin{array}{c} \log(1 + \lambda^2) \\ \lambda^2 \\ \lambda\sqrt{1 + \lambda^2} \end{array} \right]^T \boldsymbol{\eta}_\lambda \right\} \quad (16)$$

for some 3×1 vector $\boldsymbol{\eta}_\lambda$. Priors of the form $\lambda \sim N(0, \sigma_\lambda^2)$ are allowable under conjugacy constraints since these are a special case of (16) with $\boldsymbol{\eta}_\lambda = [0 \quad -1/(2\sigma_\lambda^2) \quad 0]^T$.

The message natural parameter updates depend on the first derivative of

$$\zeta(x) \equiv \log\{2\Phi(x)\} \quad \text{which leads to} \quad \zeta'(x) \equiv \frac{\phi(x)}{\Phi(x)}.$$

Software such as the function `zeta()` within the package `sn` (Azzalini, 2017) of the R computing environment (R Core Team, 2017) supports stable computation of ζ' .

Algorithm 4 provides the inputs, updates and outputs for the Skew Normal likelihood fragments. The $(\mathbf{ET})_2^{\text{SS}}$ and $(\mathbf{ET})_3^{\text{SS}}$ notation is explained in Section S.2.5 of the online supplement.

Justification for Algorithm 4 is given in Section S.3.4 of the online supplement.

3.5 Finite Normal Mixture Likelihood

The Finite Normal Mixture likelihood fragments involve the likelihood

$$y_i | \boldsymbol{\theta}, \sigma^2 \stackrel{\text{ind.}}{\sim} \text{Normal-Mixture}\left((\mathbf{A}\boldsymbol{\theta})_i, \sigma, \mathbf{w}, \mathbf{m}, \mathbf{s}\right), \quad 1 \leq i \leq n, \quad (17)$$

where \mathbf{w}, \mathbf{m} and \mathbf{s} are each constant $K \times 1$ vectors. Finite Normal Mixture approximation of difficult response density functions can be a “last resort” for development of tractable Bayesian inference algorithms. See, for example, Frühwirth-Schnatter and Wagner (2006) and Frühwirth-Schnatter et al. (2009). In the variational inference context, Wand et al. (2011) showed how Finite Normal Mixture approximation benefits variational inference for Generalized Extreme Value response models.

If we introduce auxiliary random variables $\mathbf{a}_i \equiv (a_{i1}, \dots, a_{iK})^T$ such that

$$\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \text{Multinomial}(1, \mathbf{w}), \quad 1 \leq i \leq n,$$

then we can re-express (17) as

$$p(\mathbf{y} | \boldsymbol{\theta}, \sigma^2, \mathbf{a}_1, \dots, \mathbf{a}_n) = \prod_{i=1}^n \prod_{k=1}^K \left[\sigma^{-1} (2\pi s_k^2)^{-1/2} \exp \left\{ -\frac{1}{2s_k^2} \left(\frac{(\mathbf{y} - \mathbf{A}\boldsymbol{\theta})_i}{\sigma} - m_k \right)^2 \right\} \right]^{a_{ik}}, \quad (18)$$

$$\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \text{Multinomial}(1, \mathbf{w}).$$

Even though the \mathbf{a}_i are vectors, we will use the abbreviation $\mathbf{a} \equiv \mathbf{a}_1, \dots, \mathbf{a}_n$ from now onwards. The q -density product form we consider is

$$q(\boldsymbol{\theta}, \sigma^2, \mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2)q(\mathbf{a}) = q(\boldsymbol{\theta})q(\sigma^2) \prod_{i=1}^n q(\mathbf{a}_i).$$

Algorithm 4 The inputs, updates and outputs of the Skew Normal likelihood fragments.

Data Inputs: \mathbf{y}, \mathbf{A} .

Parameter Inputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}), \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \sigma^2,$
 $\boldsymbol{\eta}_{\sigma^2} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \lambda, \sigma^2, \mathbf{a}), \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \lambda, \boldsymbol{\eta}_{\lambda} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a}).$

Updates:

$$\begin{aligned} \mu_{q(1/\sigma)} &\leftarrow (ET)_2^{\text{ISRN}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2) \\ \mu_{q(1/\sigma^2)} &\leftarrow (ET)_3^{\text{ISRN}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2) \\ \mu_{q(\lambda^2)} &\leftarrow (ET)_2^{\text{SS}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \leftrightarrow \lambda) \\ \mu_{q(\lambda\sqrt{1+\lambda^2})} &\leftarrow (ET)_3^{\text{SS}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \leftrightarrow \lambda) \\ \boldsymbol{\omega}_{10} &\leftarrow \mathbf{y} + \frac{1}{2}\mathbf{A}\left\{\text{vec}^{-1}\left(\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}\right)_2\right)\right\}^{-1}\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}\right)_1 \\ \boldsymbol{\omega}_{11} &\leftarrow G_{\text{VMP}}\left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \mathbf{A}, \mathbf{A}^T \mathbf{y}, \mathbf{y}^T \mathbf{y}\right) \\ \boldsymbol{\omega}_{12} &\leftarrow \frac{\mu_{q(1/\sigma)} \mu_{q(\lambda\sqrt{1+\lambda^2})} \boldsymbol{\omega}_{10}}{\sqrt{1 + \mu_{q(\lambda^2)}}}; \quad \boldsymbol{\omega}_{13} \leftarrow \frac{\boldsymbol{\omega}_{12} + \zeta'(\boldsymbol{\omega}_{12})}{\sqrt{1 + \mu_{q(\lambda^2)}}} \\ \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \boldsymbol{\theta} &\leftarrow \{1 + \mu_{q(\lambda^2)}\} \mu_{q(1/\sigma^2)} \begin{bmatrix} \mathbf{A}^T \mathbf{y} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \mathbf{A}) \end{bmatrix} \\ &\quad - \mu_{q(\lambda\sqrt{\lambda^2+1})} \mu_{q(1/\sigma)} \begin{bmatrix} \mathbf{A}^T \boldsymbol{\omega}_{13} \\ \mathbf{0} \end{bmatrix} \\ \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \sigma^2 &\leftarrow \begin{bmatrix} -n/2 \\ \mu_{q(\lambda\sqrt{1+\lambda^2})} \boldsymbol{\omega}_{10}^T \boldsymbol{\omega}_{13} \\ \{1 + \mu_{q(\lambda^2)}\} \boldsymbol{\omega}_{11} \end{bmatrix} \\ \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \lambda &\leftarrow \begin{bmatrix} n/2 \\ \mu_{q(1/\sigma^2)} \boldsymbol{\omega}_{11} - \frac{n + \mathbf{1}_n^T [\boldsymbol{\omega}_{12} \odot \{\boldsymbol{\omega}_{12} + \zeta'(\boldsymbol{\omega}_{12})\}]}{2\{1 + \mu_{q(\lambda^2)}\}} \\ \mu_{q(1/\sigma)} \boldsymbol{\omega}_{10}^T \boldsymbol{\omega}_{13} \end{bmatrix} \end{aligned}$$

Parameter Outputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \sigma^2, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \lambda, \mathbf{a})} \rightarrow \lambda$

The factor graph fragments for the Finite Normal Mixture likelihood are shown in Figure 6.

As in Sections 3.3 and 3.4, the conjugate distribution for σ^2 is the Inverse Square Root Nadarajah distribution (Section S.2.3 of the online supplement).

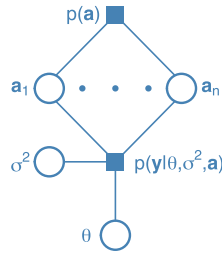


Figure 6: Fragments for the Finite Normal Mixture likelihood specification with independent Multinomial(1, \mathbf{w}) auxiliary variables $\mathbf{a}_1, \dots, \mathbf{a}_n$.

The inputs, updates and outputs for the Finite Normal Mixture likelihood fragments are listed in Algorithm 5, and justifications are in Section S.3.5 of the online supplement.

3.6 Support Vector Machine Pseudo-likelihood

Luts and Ormerod (2014) derived mean field variational Bayes algorithms for support vector machine classification using the auxiliary variable representation of the hinge loss psuedo-likelihood of Polson and Scott (2011). The approach is founded upon the following result:

$$\int_0^\infty (2\pi a)^{-1/2} \exp\left\{-\frac{(1+a-x)^2}{2a}\right\} da = \exp\{-2(1-x)_+\}, \tag{19}$$

where $u_+ \equiv \max(0, u)$ for any $u \in \mathbb{R}$. Letting $I(\mathcal{P})$ be the indicator of whether the proposition \mathcal{P} is true, note that (19) can be re-expressed as follows:

$$\begin{aligned} \text{if } p(x|a) \text{ is the } N(a+1, a) \text{ density function in } x \text{ and } \check{p}(a) = I(a > 0) \text{ then} \\ \check{p}(x) \equiv \int_{-\infty}^\infty p(x|a)\check{p}(a) da = \exp\{-2(1-x)_+\}. \end{aligned} \tag{20}$$

In (20) the pseudo-density function $\check{p}(x)$ is represented as a mixture of a particular Normal density function and the auxiliary variable pseudo-density function $\check{p}(a)$. As we will see, such a representation is amenable to the VMP updating equations with pseudo-density functions treated as ordinary density functions. As explained in Polson and Scott (2011), the hinge loss pseudo-density function could be replaced by an ordinary density function via normalization. However, the pseudo-density function version leads to the traditional support vector machine classifier.

The Support Vector Machine pseudo-likelihood fragments are concerned with the pseudo-likelihood specification

$$\check{p}(\mathbf{y}|\boldsymbol{\theta}) = \prod_{i=1}^n \exp[-2\{1 - (2y_i - 1)(\mathbf{A}\boldsymbol{\theta})_i\}_+], \tag{21}$$

Algorithm 5 The inputs, updates and outputs of the Finite Normal Mixture likelihood fragments.

Data Inputs: $\mathbf{y}, \mathbf{A}, K, w, m, s$.

Parameter Inputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a}), \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2,$
 $\boldsymbol{\eta}_{\sigma^2} \rightarrow p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})$.

Updates:

$$\mu_{q(1/\sigma)} \leftarrow (ET)_2^{\text{ISRN}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2); \mu_{q(1/\sigma^2)} \leftarrow (ET)_3^{\text{ISRN}}(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \sigma^2)$$

$$\boldsymbol{\omega}_{15} \leftarrow \mathbf{y} + \frac{1}{2} \mathbf{A} \left\{ \text{vec}^{-1} \left((\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta})_2 \right) \right\}^{-1} (\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta})_1$$

$$\boldsymbol{\omega}_{16} \leftarrow \left[G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{A}, \mathbf{A}^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{y}, y_i^2 \right) \right]_{1 \leq i \leq n}$$

$$\boldsymbol{\Omega}_{17} \leftarrow \mu_{q(1/\sigma)} \boldsymbol{\omega}_{15} (\mathbf{m}/s^2)^T + \mu_{q(1/\sigma^2)} \boldsymbol{\omega}_{16} (\mathbf{1}_K/s^2)^T + \mathbf{1}_n \{ \log(w/s) - (m^2)/(2s^2) \}^T$$

$$\boldsymbol{\Omega}_{18} \leftarrow \exp(\boldsymbol{\Omega}_{17}) / \{ \exp(\boldsymbol{\Omega}_{17}) \mathbf{1}_K \mathbf{1}_K^T \}; \quad \boldsymbol{\omega}_{19} \leftarrow \boldsymbol{\Omega}_{18} (\mathbf{1}_K/s^2)$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta} \leftarrow \mu_{q(1/\sigma^2)} \begin{bmatrix} \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_{19}) \mathbf{y} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_{19}) \mathbf{A}) \end{bmatrix} - \mu_{q(1/\sigma)} \begin{bmatrix} \mathbf{A}^T \boldsymbol{\Omega}_{18} (\mathbf{m}/s^2) \\ \mathbf{0} \end{bmatrix}$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2 \leftarrow \begin{bmatrix} -n/2 \\ \boldsymbol{\omega}_{15}^T \boldsymbol{\Omega}_{18} (\mathbf{m}/s^2) \\ G_{\text{VMP}} \left(\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \leftrightarrow \boldsymbol{\theta}; \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_{19}) \mathbf{A}, \right. \\ \left. \mathbf{A}^T \text{diag}(\boldsymbol{\omega}_{19}) \mathbf{y}, \mathbf{y}^T \text{diag}(\boldsymbol{\omega}_{19}) \mathbf{y} \right) \end{bmatrix}$$

Parameter Outputs: $\boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{p(\mathbf{y}|\boldsymbol{\theta}, \sigma^2, \mathbf{a})} \rightarrow \sigma^2$.

where the $y_i \in \{0, 1\}$ are indicators of class membership in a two-class classification setting. If we now introduce an auxiliary variable vector $\mathbf{a} = (a_1, \dots, a_n)$ with entries a_i , $1 \leq i \leq n$, with each independently having the pseudo-density function $\check{p}(a_i) = I(a_i > 0)$ then, using (20), (21) is equivalent to

$$\check{p}(\mathbf{y}|\mathbf{a}, \boldsymbol{\theta}) = \prod_{i=1}^n (2\pi a_i)^{-1/2} \exp \left[-\frac{\{1 + a_i - (2y_i - 1)(\mathbf{A}\boldsymbol{\theta})_i\}^2}{2a_i} \right], \quad (22)$$

$$\check{p}(\mathbf{a}) = \prod_{i=1}^n I(a_i > 0).$$

The corresponding factor graph fragments are shown in Figure 7.

Under the assumption that all messages passed to $\boldsymbol{\theta}$ are Multivariate Normal, Algorithm 6 provides updates for the natural parameter vector passed from $\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})$

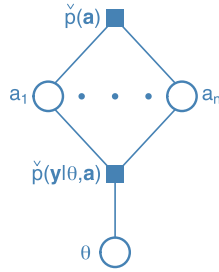


Figure 7: Fragments for the Support Vector Machine pseudo-likelihood specification with independent auxiliary variables $\mathbf{a} = (a_1, \dots, a_n)$ having psuedo-density function $\check{p}(\mathbf{a}) = \prod_{i=1}^n I(a_i > 0)$.

Algorithm 6 The inputs, updates and outputs of the Support Vector Machine pseudo-likelihood fragments.

Data Inputs: \mathbf{y}, \mathbf{A} .

Parameter Inputs: $\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \rightarrow \boldsymbol{\theta}, \boldsymbol{\eta}_{\boldsymbol{\theta}} \rightarrow \check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})$

Updates:

$$\boldsymbol{\omega}_{20} \leftarrow -\frac{1}{2} \mathbf{A} \left\{ \text{vec}^{-1} \left(\left(\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \leftrightarrow \boldsymbol{\theta} \right)_2 \right) \right\}^{-1} \left(\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \leftrightarrow \boldsymbol{\theta} \right)_1 \Big]$$

$$\boldsymbol{\omega}_{21} \leftarrow -\frac{1}{2} \text{diagonal} \left[\mathbf{A} \left\{ \text{vec}^{-1} \left(\left(\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \leftrightarrow \boldsymbol{\theta} \right)_2 \right) \right\}^{-1} \mathbf{A}^T \right]$$

$$\boldsymbol{\omega}_{22} \leftarrow \left[\{(2\mathbf{y} - \mathbf{1}_n) \odot \boldsymbol{\omega}_{20} - \mathbf{1}_n\}^2 + \boldsymbol{\omega}_{21} \right]^{-1/2}$$

$$\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \rightarrow \boldsymbol{\theta} \leftarrow \begin{bmatrix} \mathbf{A}^T \{(\mathbf{1}_n + \boldsymbol{\omega}_{22}) \odot (2\mathbf{y} - \mathbf{1})\} \\ -\frac{1}{2} \text{vec}(\mathbf{A}^T \text{diag}(\boldsymbol{\omega}_{22}) \mathbf{A}) \end{bmatrix}$$

Parameter Outputs: $\boldsymbol{\eta}_{\check{p}(\mathbf{y}|\boldsymbol{\theta}, \mathbf{a})} \rightarrow \boldsymbol{\theta}$

to $\boldsymbol{\theta}$. An attractive feature of the Support Vector Machine pseudo-likelihood fragment updates is that each of them are simple closed form operations.

4 Illustrations

We now provide some illustrations of how the fragment updates of Section 3 can be used to move from one variational inference analysis to another, without having to start from scratch.

4.1 Ordinary to Quantile Nonparametric Regression

First consider ordinary nonparametric regression via the Bayesian mixed model-based penalized spline model used in Section 3.2.1 of Wand (2017). We quickly recap the details here. The data are the predictor/response pairs (x_i, y_i) , $1 \leq i \leq n$, and the nonparametric regression model is:

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

where the model for the mean function f takes the form

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K u_k z_k(x) \quad \text{with} \quad u_k | \sigma_u^2 \stackrel{\text{ind.}}{\sim} N(0, \sigma_u^2) \quad (23)$$

and $\{z_k : 1 \leq k \leq K\}$ is a suitable spline basis. The full model used in Wand (2017) is

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\varepsilon^2 &\sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\varepsilon^2 \mathbf{I}), \quad \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} | \sigma_u^2 \sim N\left(\begin{bmatrix} \boldsymbol{\mu}_\beta \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_\beta & \mathbf{0} \\ \mathbf{0} & \sigma_u^2 \mathbf{I} \end{bmatrix}\right), \\ \sigma_u^2 | a_u &\sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A_u^2), \\ \sigma_\varepsilon^2 | a_\varepsilon &\sim \text{Inverse-}\chi^2(1, 1/a_\varepsilon), \quad a_\varepsilon \sim \text{Inverse-}\chi^2(1, 1/A_\varepsilon^2) \end{aligned} \quad (24)$$

where

$$\mathbf{X} \equiv \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \text{and} \quad \mathbf{Z} \equiv \begin{bmatrix} z_1(x_1) & \cdots & z_K(x_1) \\ \vdots & \ddots & \vdots \\ z_1(x_n) & \cdots & z_K(x_n) \end{bmatrix}.$$

The 2×1 vector $\boldsymbol{\mu}_\beta$, 2×2 symmetric positive definite matrix $\boldsymbol{\Sigma}_\beta$ and the positive numbers A_u and A_ε are user-specified hyperparameters. Note that

$$\sigma_u^2 | a_u \sim \text{Inverse-}\chi^2(1, 1/a_u), \quad a_u \sim \text{Inverse-}\chi^2(1, 1/A_u^2)$$

is equivalent to σ_u having a Half Cauchy prior with scale parameter A_u , but this auxiliary variable representation has advantages for VMP fitting. The final choice is the form of the z_k and the value of K . In the upcoming example we used canonical cubic O'Sullivan splines (Wand and Ormerod, 2008) with $K = 27$.

The joint posterior density function is approximated according to the following product restriction

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma_u^2, a_u, \sigma_\varepsilon^2, a_\varepsilon | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}) q(\sigma_u^2) q(a_u) q(\sigma_\varepsilon^2) q(a_\varepsilon). \quad (25)$$

VMP fitting of (24) can be accomplished by using the natural parameter updates for each of the fragments described in Section 4.1 of Wand (2017). The relevant factor graph is in the left panel of Figure 8 with the Gaussian likelihood fragment shown in red. We applied the VMP fitting procedure to data on 4,847 Zambian children from a

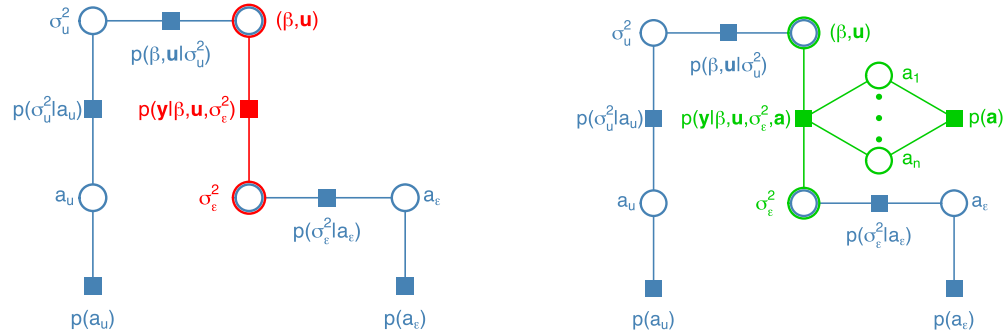


Figure 8: Left panel: Factor graph for the ordinary nonparametric regression model. The Gaussian likelihood fragment is shown in red. Right panel: Factor graph for the quantile nonparametric regression model. The Asymmetric Laplace likelihood fragments are shown in green.

1992 demographic and health survey. These data are part of the data frame `Zambia` in the R package `INLA` (Rue et al., 2016). The predictor and response data are

$$\begin{aligned}
 x_i &= \text{age of the } i\text{th child in months,} \\
 \text{and } y_i &= \text{undernutrition score of the } i\text{th child, } 1 \leq i \leq 4,847.
 \end{aligned}
 \tag{26}$$

All data were standardized and the hyperparameters we set at $\mu_\beta = \mathbf{0}$, $\Sigma_\beta = 10^{10} \mathbf{I}$ and $A_u = A_\epsilon = 10^5$. The fits were back-transformed to the original units for plotting. The estimated nonparametric regression function and corresponding pointwise 95% credible set are shown in the left panel of Figure 9. The estimate shows mean undernutrition falling during the infancy period of the children before levelling off at about 2 years of age.

Now suppose that 100 τ % quantile nonparametric regression for the same data is of interest. This involves replacement of

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma_\epsilon^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma_\epsilon^2 \mathbf{I})$$

by

$$y_i | \boldsymbol{\beta}, \mathbf{u}, \sigma^2, \mathbf{a} \stackrel{\text{ind.}}{\sim} N\left((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_i + \frac{(\frac{1}{2} - \tau)\sigma}{a_i\tau(1 - \tau)}, \frac{\sigma^2}{a_i\tau(1 - \tau)}\right), \quad a_i \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(2, 1)$$

in model (24). In terms of factor graphs it involves replacement of the Gaussian likelihood fragment by the Asymmetric Laplace likelihood fragments of Figure 4. The new fragments are shown in green in the right panel of Figure 8. The VMP updates corresponding to messages away from the likelihood are identical for both models. Algorithm 3 is used for the quantile nonparametric regression fitting and inference.

As a check, the same models were fit to the data using Markov chain Monte Carlo. The nonparametric regression and quantile regression curves are very close to their

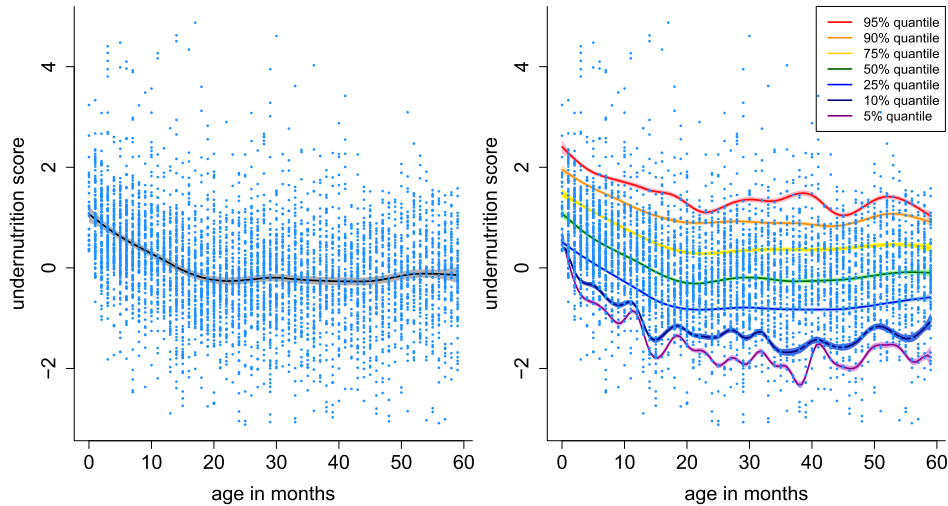


Figure 9: Left panel: VMP nonparametric regression fit to the variables on Zambian children given by (26). The curve is the approximate posterior mean and the shaded region corresponds to pointwise approximate 95% credible sets. The estimates are based on VMP applied to model (24) according to product restriction (25). The relevant factor graph is shown in the left panel of Figure 8. Right panel: VMP quantile nonparametric regression fits to the same data. The curves and shaded regions have the same definitions as for the left panel.

VMP counterparts. However, the 95% credible set bands are narrower for VMP. This is a consequence of the loss of inferential accuracy incurred by variational approximations involving auxiliary variables (see e.g. Wand et al., 2011).

4.2 Poisson to Negative Binomial Additive Model Analysis

Our second illustration involves additive model analysis when the response variable is a count. First we carried out a Poisson additive model analysis similar to those described in Section 12.3 of Ruppert et al. (2003). The data involve daily ragweed pollen counts in Kalamazoo, U.S.A., during the 1991–1994 ragweed seasons. The model is of the form

$$y_i \stackrel{\text{ind.}}{\sim} \text{Poisson}\left\{\exp\left(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + f_{z_i}(x_{4i})\right)\right\}, \quad 1 \leq i \leq n, \quad (27)$$

where $n = 334$ is the total number of days when ragweed pollen was in season during 1991–1994. The variables appearing in (27) are ragweed pollen count on the i th day (y_i), temperature residual on the i th day (x_{1i}), indicator of significant rain on the i th day (x_{2i}), wind speed in knots on the i th day (x_{3i}), day number of ragweed pollen season for the current year on which y_i was recorded (x_{4i}) and a categorical variable for the year in which y_i was recorded (one of 1991, 1992, 1993 or 1994) (z_i). Here temperature residuals are the residuals from fitting penalized splines, each having 5 effective degrees

of freedom, to temperature (in degrees Fahrenheit) versus day number for each annual ragweed pollen season. Note that (27) is not an additive model in the usual sense since $f_{z_i}(x_{4i})$ represents an interaction between year and day in ragweed pollen season. Mixed model-based penalized splines analogous to (23) are used for modelling the f_z , $z \in \{1992, 1992, 1993, 1994\}$. Let $\sigma_{u\ell}^2$, $1 \leq \ell \leq 4$, denote the variance parameters used to penalize each of the four penalized splines. The full model is

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u} \sim \text{Poisson}\{\exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})\},$$

$$\begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} \Big| \sigma_{u1}^2, \sigma_{u2}^2, \sigma_{u3}^2, \sigma_{u4}^2 \sim N \left(\begin{bmatrix} \boldsymbol{\mu}_\beta \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_\beta & \mathbf{0} \\ \mathbf{0} & \text{blockdiag}(\sigma_{u\ell}^2 \mathbf{I})_{1 \leq \ell \leq 4} \end{bmatrix} \right), \quad (28)$$

$$\sigma_{u\ell}^2 | a_{u\ell} \sim \text{Inverse-}\chi^2(1, 1/a_{u\ell}), \quad a_{u\ell} \sim \text{Inverse-}\chi^2(1, 1/A_{u\ell}^2), \quad 1 \leq \ell \leq 4.$$

Here

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{41} & I(z_1=1992) & x_{41}I(z_1=1992) & \cdots & I(z_1=1994) & x_{41}I(z_1=1994) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & \cdots & x_{4n} & I(z_n=1992) & x_{4n}I(z_n=1992) & \cdots & I(z_n=1994) & x_{4n}I(z_n=1994) \end{bmatrix}$$

and $\mathbf{Z} = [\mathbf{Z}_{1991} \mathbf{Z}_{1992} \mathbf{Z}_{1993} \mathbf{Z}_{1994}]$ where \mathbf{Z}_{1991} is an $n \times K$ matrix with (i, k) entry equal to $I(z_i = 1991)z_k(x_{4i})$ and $\mathbf{Z}_{1992}, \dots, \mathbf{Z}_{1994}$ are defined analogously. The $\boldsymbol{\beta}$ and \mathbf{u} vectors contain the coefficients to match the columns of \mathbf{X} and \mathbf{Z} respectively.

Despite the simplicity of Poisson response regression models, it is often the case that the Poisson likelihood is inadequate for modeling count response data that typically arises in practice. The crux of this inadequacy is the Poisson distribution restriction of the variance equalling the mean. It is common for the variability of count responses to be much higher than that imposed by the Poisson likelihood. If such overdispersion is ignored then standard errors are underestimated and valid statistical inference is compromised. The Negative Binomial family is an extension of the Poisson family that allows for the variance to exceed the mean. The move from this Poisson additive model to a Negative Binomial additive model involves replacement of

$$\mathbf{y} | \boldsymbol{\beta}, \mathbf{u} \sim \text{Poisson}\{\exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})\}$$

by

$$y_i | a_i \stackrel{\text{ind.}}{\sim} \text{Poisson}(a_i), \quad a_i | \boldsymbol{\beta}, \mathbf{u}, \kappa \stackrel{\text{ind.}}{\sim} \text{Gamma}[\kappa, \kappa \exp\{-(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_i\}],$$

which corresponds to the likelihood specification

$$y_i | \boldsymbol{\beta}, \mathbf{u}, \kappa \stackrel{\text{ind.}}{\sim} \text{Negative-Binomial}[\exp\{-(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_i\}, \kappa].$$

Figure 10 shows the old and the new factor graphs according to this replacement. Almost all of the fragments in these factor graphs are covered by Wand (2017) and Algorithm 1. The exception is the fragment containing the factor $p(\kappa)$, which corresponds

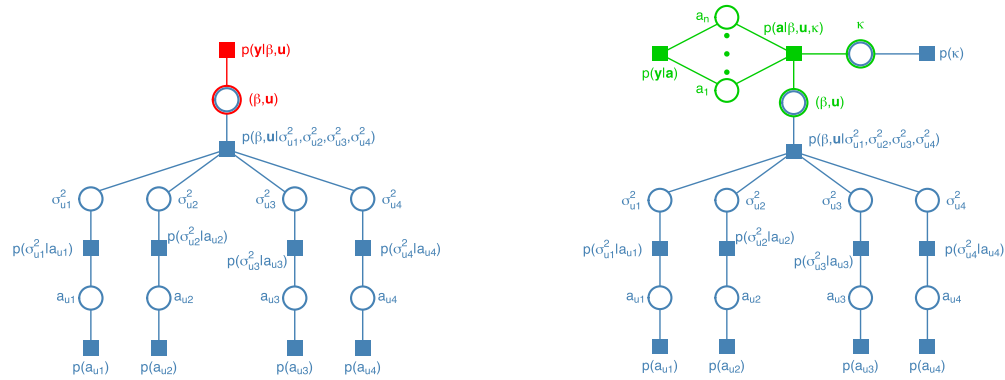


Figure 10: Left panel: Factor graph for the Poisson additive regression model. The Poisson likelihood fragment is shown in red. Right panel: Factor graph for the Negative Binomial additive model. The Negative Binomial likelihood fragments are shown in green.

to placing a prior distribution on κ . In the ragweed data analysis we used the prior $p(\kappa) = 0.01 \exp(-0.01\kappa)$, $\kappa > 0$, which implies that the message sent from $p(\kappa)$ to κ is

$$m_{p(\kappa)} \rightarrow \kappa(\kappa) = \exp \left\{ \left[\begin{array}{c} \kappa \log(\kappa) - \log\{\Gamma(\kappa)\} \\ \kappa \end{array} \right]^T \left[\begin{array}{c} 0 \\ -0.01 \end{array} \right] \right\}.$$

This prior and message simply correspond to the Exponential distribution with rate parameter 0.01. We use the Moon Rock-type representation since it is conjugate with messages passed from $p(\mathbf{a}|\boldsymbol{\beta}, \mathbf{u}, \kappa)$ to κ .

Figure 11 provides some visual summaries of the model fits. The first row shows posterior density functions for the coefficients of the predictors that enter the models linearly, and the Negative Binomial shape parameter. The posterior density functions for the Poisson model are considerably narrower than those for the Negative Binomial model, which is indicative of overdispersion being ignored in the former model. In the same vein, the posterior density function of κ has most of its support between 1.4 and 2.2. Such low κ values indicate superiority of the Negative Binomial model since the Poisson model corresponds to the $\kappa \rightarrow \infty$ limiting case.

The lower four panels of Figure 11 show the estimates of $f_{1991}, \dots, f_{1994}$ for the Poisson and Negative Binomial models. The solid curves correspond to the posterior mean for each day in season value, while the dashed curves are pointwise 95% credible sets according to the VMP approximation. The estimates are similar for each model, but the credible set bands are narrower for the Poisson model, in keeping with their ignorance of overdispersion.

Computing times for the Poisson and Negative Binomial additive models were also compared. All computing was performed using version 3.4.1 of the R language (R Core

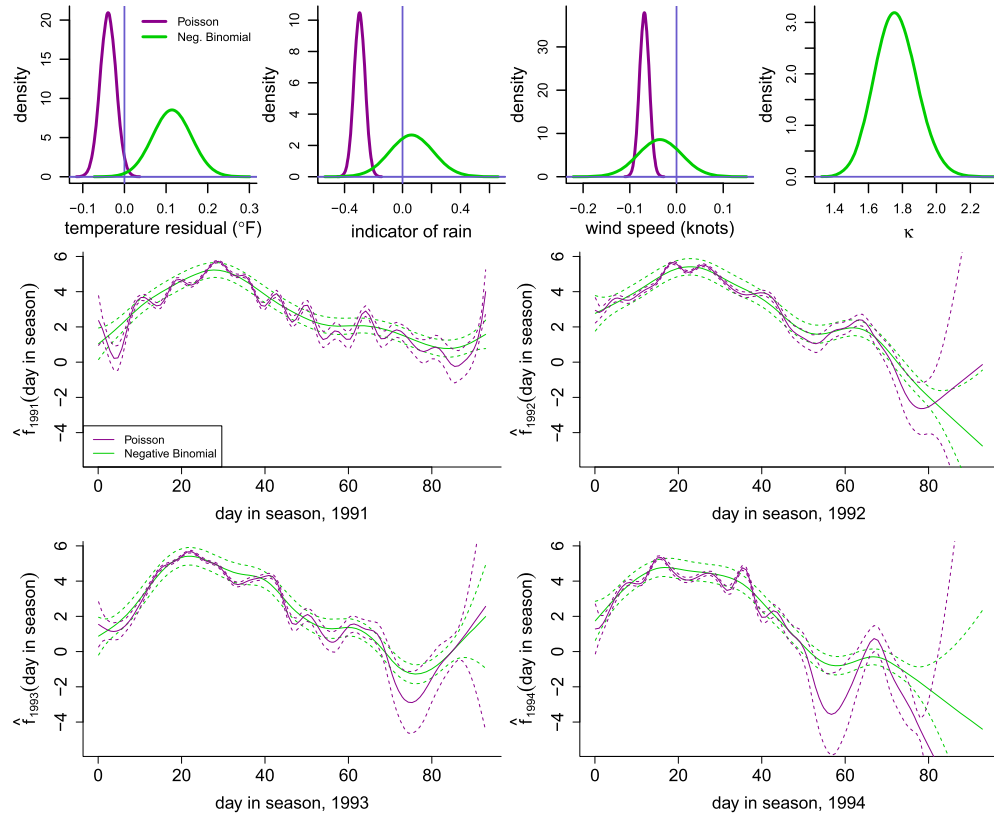


Figure 11: First three panels: VMP-approximate posterior density functions of the coefficients of temperature residual, indicator of rain and wind speed for both the Poisson additive model (28) and the Negative Binomial additive model for the ragweed data example. Fourth panel: VMP-approximate posterior density function of the κ parameter for the Negative Binomial additive model. Lower four panels: VMP-based estimates of $f_{1992}, \dots, f_{1994}$ according to each model. The solid curves are posterior means and the dashed curves are pointwise 95% credible intervals based on VMP approximate inference.

Team, 2017) on a desktop personal computer with 8 gigabytes of random access memory and a 3.2 gigahertz processor. Firstly, we determined that 250 iterations were sufficient for convergence of VMP for each model. The elapsed times were 5.5 seconds for the Poisson model and 6.9 seconds for the Negative Binomial model.

We also compared the VMP-approximate posterior density functions and additive model components with those obtained using Markov chain Monte Carlo. Excellent agreement was observed in almost all cases. An exception concerned the posterior density function for κ , with VMP under-approximating the posterior standard deviation. This phenomenon was also observed in Luts and Wand (2015).

5 Closing Remarks

As exemplified in Section 4, the algorithms presented in Section 3 concerning fragments updates for elaborate likelihoods greatly enhances the utility of VMP for semiparametric regression analyses. In addition to the primitives for VMP-based semiparametric regression laid down in Wand (2017) we have identified a small set of new primitives, corresponding to sufficient statistic expectations of the Inverse Square Root Nadarajah, Moon Rock and Sea Sponge distributions. Once their computation is established in a suite of computer programmes, a much richer class of models can be handled via the VMP paradigm.

Supplementary Material

Supplement for: Variational Message Passing for Elaborate Response Regression Models (DOI: [10.1214/18-BA1098SUPP](https://doi.org/10.1214/18-BA1098SUPP); .pdf).

References

- Azzalini, A. (2017). *The R package sn: The skew-normal and related distributions, such as the skew-t (version 1.5)*. URL <http://azzalini.stat.unipd.it/SN> MR3468021. 385
- Azzalini, A. and Dalla Valle, A. (1996). “The multivariate skew-normal distribution.” *Biometrika*, 83: 715–726. MR1440039. doi: <https://doi.org/10.1093/biomet/83.4.715>. 383
- Frühwirth-Schnatter, S., Frühwirth, R., Held, L., and Rue, H. (2009). “Improved auxiliary mixture sampling for hierarchical models of non-Gaussian data.” *Statistics and Computing*, 19: 479–492. MR2565319. doi: <https://doi.org/10.1007/s11222-008-9109-4>. 385
- Frühwirth-Schnatter, S. and Pyne, S. (2010). “Bayesian inference for finite mixtures of univariate and multivariate skew-normal and skew-t distributions.” *Biostatistics*, 11: 317–336. 382
- Frühwirth-Schnatter, S. and Wagner, H. (2006). “Auxiliary mixture sampling for parameter-driven models of time series of counts with applications to state space modelling.” *Biometrika*, 93: 827–841. MR2285074. doi: <https://doi.org/10.1093/biomet/93.4.827>. 385
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. W. (2013). “Stochastic variational inference.” *Journal of Machine Learning Research*, 14: 1303–1347. MR3081926. 372
- Knowles, D. A. and Minka, T. (2011). “Non-conjugate variational message passing for multinomial and binary regression.” In *Advances in Neural Information Processing Systems*, 1701–1709. 377

- Kotz, S., Kozubowski, T. J., and Podgórski, K. (2001). *The Laplace Distribution and Generalizations*. Boston: Birkhäuser. MR1935481. doi: <https://doi.org/10.1007/978-1-4612-0173-1>. 380
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. (2017). “Automatic differentiation variational inference.” *Journal of Machine Learning Research*, 18: 1–45. MR3634881. 372
- Lachos, V. H., Ghosh, P., and Arellano-Valle, R. B. (2010). “Likelihood based inference for skew-normal independent linear mixed models.” *Statistica Sinica*, 303–322. MR2640696. 382
- Lange, K. L., Little, R. J. A., and Taylor, J. M. G. (1989). “Robust statistical modeling using the t distribution.” *Journal of the American Statistical Association*, 84: 881–896. MR1134486. 379
- Luts, J. and Ormerod, J. T. (2014). “Mean field variational Bayesian inference for support vector machine classification.” *Computational Statistics & Data Analysis*, 73: 163–176. MR3147981. doi: <https://doi.org/10.1016/j.csda.2013.10.030>. 387
- Luts, J. and Wand, M. P. (2015). “Variational inference for count response semiparametric regression.” *Bayesian Analysis*, 10: 991–1023. MR3432247. doi: <https://doi.org/10.1214/14-BA932>. 377, 395
- McLean, M. W. and Wand, M. P. (2018). “Supplement for: Variational Message Passing for Elaborate Response Regression Models.” *Bayesian Analysis*. doi: <https://doi.org/10.1214/18-BA1098SUPP>. 372
- Minka, T. (2005). “Divergence measures and message passing.” *Microsoft Research Technical Report Series*, MSR-TR-2005-173: 1–17. 371
- Minka, T. and Winn, J. (2008). “Gates: A graphical notation for mixture models.” *Microsoft Research Technical Report Series*, MSR-TR-2008-185: 1–16. 371
- Nadarajah, S. (2008). “A new model for symmetric and skewed data.” *Probability in the Engineering and Informational Sciences*, 22: 261–271. MR2399727. doi: <https://doi.org/10.1017/S0269964808000156>. 382
- Ormerod, J. T. and Wand, M. P. (2010). “Explaining variational approximations.” *The American Statistician*, 64: 140–153. MR2757005. doi: <https://doi.org/10.1198/tast.2010.09058>. 373
- Polson, N. G. and Scott, S. L. (2011). “Data augmentation for support vector machines.” *Bayesian Analysis*, 6: 1–23. MR2781803. doi: <https://doi.org/10.1214/11-BA601>. 387
- R Core Team (2017). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/> 385, 394
- Rue, H., Martino, S., Lindgren, F., Simpson, D., Riebler, A., and Krainski, E. (2016). *The R package ‘INLA’: Functions which allow to perform full Bayesian analysis of*

- latent Gaussian models using integrated nested Laplace approximation (version 0.0)*.
URL <http://www.r-inla.org> 391
- Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric Regression*. New York: Cambridge University Press. MR1998720. doi: <https://doi.org/10.1017/CB09780511755453>. 372, 392
- Tipping, M. E. and Lawrence, N. D. (2003). “A variational approach to robust Bayesian interpolation.” In *Institute of Electrical and Electronics Engineers Workshop of Neural Networks for Signal Processing*, 229–238. 379
- Titsias, M. K. and Lázaro-Gredilla, M. (2014). “Doubly stochastic variational Bayes for non-conjugate inference.” *Proceedings of Machine Learning Research*, 32: 1971–1979. 372
- Verdinelli, I. and Wasserman, L. (1991). “Bayesian analysis of outlier problems using the Gibbs sampler.” *Statistics and Computing*, 1: 105–117. 379
- Wand, M. P. (2017). “Fast approximate inference for arbitrarily large semiparametric regression models via message passing (with discussion).” *Journal of the American Statistical Association*, 112: 137–168. MR3646558. doi: <https://doi.org/10.1080/01621459.2016.1197833>. 371, 372, 373, 374, 375, 390, 393, 396
- Wand, M. P. and Ormerod, J. T. (2008). “On semiparametric regression with O’Sullivan penalized splines.” *Australian & New Zealand Journal of Statistics*, 50: 179–198. MR2431193. doi: <https://doi.org/10.1111/j.1467-842X.2008.00507.x>. 390
- Wand, M. P., Ormerod, J. T., Padoan, S. A., and Frühwirth, R. F. (2011). “Mean field variational Bayes for elaborate distributions.” *Bayesian Analysis*, 6: 847–900. MR2869967. doi: <https://doi.org/10.1214/11-BA631>. 385, 392
- Winn, J. and Bishop, C. M. (2005). “Variational message passing.” *Journal of Machine Learning Research*, 6: 661–694. MR2249835. 371
- Yang, Y., Wang, H. J., and He, X. (2016). “Posterior inference in Bayesian quantile regression with Asymmetric Laplace likelihood.” *International Statistical Review*, 84: 327–344. MR3580414. doi: <https://doi.org/10.1111/insr.12114>. 380
- Yu, K. and Moyeed, R. A. (2001). “Bayesian quantile regression.” *Statistics and Probability Letters*, 54: 437–447. MR1861390. doi: [https://doi.org/10.1016/S0167-7152\(01\)00124-9](https://doi.org/10.1016/S0167-7152(01)00124-9). 380

Acknowledgments

This research was partially supported by the Australian Research Council Discovery Project DP140100441. We are grateful to Peter Alspach, Chris Jones, Luca Maestrini, John Maindonald, Saralees Nadarajah and Thomas Yee for their assistance with this research.