



Factor graph fragmentation of expectation propagation

Wilson Y. Chen¹ · Matt P. Wand¹ 

Received: 15 July 2019 / Accepted: 23 October 2019 / Published online: 1 January 2020
© Korean Statistical Society 2020

Abstract

Expectation propagation is a general approach to fast approximate inference for graphical models. The existing literature treats models separately when it comes to deriving and coding expectation propagation inference algorithms. This comes at the cost of similar, long-winded algebraic steps being repeated and slowing down algorithmic development. We demonstrate how *factor graph fragmentation* can overcome this impediment. This involves adoption of the message passing on a factor graph approach to expectation propagation and identification of factor graph sub-graphs, which we call *fragments*, that are common to wide classes of models. Key fragments and their corresponding messages are catalogued which means that their algebra does not need to be repeated. This allows compartmentalization of coding and efficient software development.

Keywords Approximate Bayesian inference · Generalized linear mixed models · Graphical models · Kullback–Leibler projection · Message passing

1 Introduction

Expectation propagation (e.g. Minka 2005) is gaining popularity as a general approach to fitting and inference for large graphical models, including those that arise in statistical contexts such as Bayesian generalized linear mixed models (e.g. Gelman et al. 2014; Kim and Wand 2018). Compared with Markov chain Monte Carlo approaches, expectation propagation has the attractions of speed and parallelizability of the computing across multiple processors making it more amenable to high volume/velocity data

Electronic supplementary material The online version of this article (<https://doi.org/10.1007/s42952-019-00033-9>) contains supplementary material, which is available to authorized users.

✉ Matt P. Wand
matt.wand@uts.edu.au

Wilson Y. Chen
wilson@ism.ac.jp

¹ University of Technology Sydney, Ultimo, Australia

applications. One price to be paid is inferential accuracy since expectation propagation uses product density simplifications of joint posterior density functions. Another is algebraic overhead: as demonstrated by Kim and Wand (2016) several pages of algebra are required to derive explicit programmable expectation propagation algorithms for even very simple Bayesian models. This article alleviates the latter cost. Using the notions of *message passing* and *factor graph fragments* we demonstrate the compartmentalization of expectation propagation algebra and coding. The resultant infrastructure and updating formulae lead to much more efficient expectation propagation fitting and inference and allows extension to arbitrarily large Bayesian models.

Expectation propagation and *mean field variational Bayes* are the two most common paradigms for obtaining fast approximate inference algorithms for graphical models (e.g. Bishop 2006; Wainwright and Jordan 2008; Murphy 2012). Each is driven by minimum Kullback–Leibler divergence considerations. As explained in Minka (2005), they can both be expressed as message passing algorithms on factor graphs. The alternative appellation *variational message passing* is used for mean field variational Bayes when such an approach is used. The software platform Infer.NET (Minka et al. 2014) uses both expectation propagation and variational message passing to perform fast approximate inference for graphical models.

The term “expectation propagation” was introduced in Minka (2001), although earlier references which contain related ideas include Thouless et al. (1977) and Opper and Winther (2000). Heskes and Zoeter (2002), Opper and Winther (2005) and Heskes et al. (2005) delved into the optimization theory aspects of expectation propagation and related it to the methodology of Thouless et al. (1977) and to the notions of *expectation consistent approximation* and *Bethe free energy*. In these articles, both *single-loop* and *double-loop* versions of expectation propagation are studied, with the latter being a more involved route towards achieving convergence guarantees. Convergence issues aside, there are two main ways by which expectation propagation is defined for approximate inference in Bayesian statistical models. The more prominent version of expectation propagation is defined for models in which the joint density function of the parameters is a product over so-called *sites*, often corresponding to data points, and *cavity* and *tilted* distributions are obtained iteratively using particular Kullback–Leibler projections (e.g. Gelman et al. 2014, Section 13.8). There may also be the restriction that the approximate posterior density function is Gaussian (e.g. Jylänki et al. 2011). The other way to define expectation propagation is via *message passing* on a *factor graph* (e.g. Minka 2005). The factor graph definition is more general than the product over sites definition. Factorizations with respect to the data are not necessary and any exponential family can arise in expectation propagation-approximate posterior density functions. A simple instructive example involves a $N(\mu, \sigma^2)$ random sample Bayesian model. As shown in Kim and Wand (2016), the message passing on factor graph approach can lead to a non-Gaussian approximate posterior density function for σ^2 and without the need for a product over sites. In Jylänki et al. (2011) the analogous variance parameter is, instead, estimated via maximization of an approximate marginal log-likelihood. On account of its generality and fragmentation advantages, in this article we use the message passing on factor graph definition of expectation propagation. Details are given in Sect. 2.4.

Recently Wand (2017) introduced factor graph fragmentation to streamline variational message passing for semiparametric regression analysis. Semiparametric regression (e.g. Ruppert et al. 2009) is a big class of flexible regression models that includes generalized linear mixed models, generalized additive models and varying-coefficient models as special cases. Nolan and Wand (2017) and McLean and Wand (2018) built on Wand (2017) for more elaborate likelihood fragments.

The crux of this article is to show how the factor graph fragment idea also can be used to streamline expectation propagation. We focus on semiparametric regression models. However, the approach is quite general and applies to other graphical models for which expectation propagation is feasible. The fragment updating algorithms presented and derived here cover a wide range of semiparametric models and pave the way for future derivations of the same type.

Section 2 provides the background material needed for factor graph fragmentation of expectation propagation. This includes exponential family and Kullback–Leibler projection theory, as well as the notions of factor graphs and their fragment sub-graphs. The article’s centerpiece is Sect. 3 in which several key fragments are identified and have message updates derived and catalogued. Such cataloguing implies that updates for a particular fragment never have to be derived again and only need to be implemented once in an expectation propagation software suite. An illustration involving generalized additive mixed model analysis of data from a longitudinal public health study is provided in Sect. 4. Section 5 contains some commentary of fragmentation of expectation propagation for more elaborate models.

2 Background material

Factor graph fragmentation of expectation propagation relies on definitions and results concerning both distribution theory and graph theory, not all of which are commonplace in the statistics literature. We provide the necessary background material in this section.

2.1 Exponential family distributions

A $d \times 1$ random vector \mathbf{x} has an exponential family distribution if its probability mass function or density function admits the form

$$p(\mathbf{x}) = \exp\{\mathbf{T}(\mathbf{x})^T \boldsymbol{\eta} - A(\boldsymbol{\eta})\}h(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \boldsymbol{\eta} \in H.$$

The vectors $\mathbf{T}(\mathbf{x})$ and $\boldsymbol{\eta}$ are called, respectively, the *sufficient statistic* and *natural parameter*. The set H is the space of allowable natural parameter values. The function $A(\boldsymbol{\eta})$ is called the *log-partition function* and $h(\mathbf{x})$ is the *base measure*. A key exponential family distributional result is that

$$E\{\mathbf{T}(\mathbf{x})\} = \nabla A(\boldsymbol{\eta}) \tag{1}$$

where $\nabla A(\boldsymbol{\eta})$ is the column vector of partial derivatives of $A(\boldsymbol{\eta})$ with respect to each of the components of $\boldsymbol{\eta}$.

Table 1 Sufficient statistics, log-partition functions, base measures and natural parameter spaces of three univariate exponential families

name	$T(x)$	$A(\eta)$	$h(x)$	H
Normal	$\begin{bmatrix} x \\ x^2 \end{bmatrix}$	$-\frac{1}{4}(\eta_1^2/\eta_2)$ $-\frac{1}{2} \log(-2\eta_2)$	1	$\{(\eta_1, \eta_2) : \eta_1 \in \mathbb{R}, \eta_2 < 0\}$
Inverse Chi-Squared	$\begin{bmatrix} \log(x) \\ 1/x \end{bmatrix}$	$(\eta_1 + 1) \log(-\eta_2)$ $+\log \Gamma(-\eta_1 - 1)$	$I(x > 0)$	$\{(\eta_1, \eta_2) : \eta_1 < -1, \eta_2 < 0\}$
Moon Rock	$\begin{bmatrix} x \log(x) \\ -\log \Gamma(x) \\ x \end{bmatrix}$	$\log \left[\int_0^\infty \{t^t / \Gamma(t)\}^{\eta_1} \times \exp(\eta_2 t) dt \right]$	$I(x > 0)$	$\{(\eta_1, \eta_2) : \eta_1 > 0, \eta_1 + \eta_2 < 0\}$

Table 1 lists each of the univariate exponential families distributions arising in this article, along with their defining functions and parameter spaces. The Normal and Inverse Chi-Squared exponential families are well known. The Moon Rock exponential family is less established, and is given this name in McLean and Wand (2018). In Table 1 and elsewhere we use the following indicator function notation: $I(\mathcal{P}) = 1$ if the proposition \mathcal{P} is true and $I(\mathcal{P}) = 0$ if \mathcal{P} is false.

Note also that the Inverse Chi-Squared exponential family is equivalent to the *Inverse Gamma* exponential family. The two families differ in their common parametrizations as explained in, for example, Section S.1.3 of the online supplement of Wand (2017). The Inverse Chi-Squared distribution has the advantage of being the special case of the Inverse Wishart distribution for 1×1 random matrices. Throughout this article we write

$$X \sim \text{Inverse-Wishart}(\kappa, \Lambda)$$

to denote a $d \times d$ random matrix X having density function

$$p(X) = C_{d,\kappa}^{-1} |\Lambda|^{\kappa/2} |X|^{-(\kappa+d+1)/2} \exp\{-\frac{1}{2}\text{tr}(\Lambda X^{-1})\} I(X \text{ symmetric and positive definite})$$

where $\kappa > d - 1$, Λ is a $d \times d$ symmetric positive definite matrix and

$$C_{d,\kappa} \equiv 2^{d\kappa/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma\left(\frac{\kappa + 1 - j}{2}\right). \tag{2}$$

For the special case of $d = 1$ we write

$$x \sim \text{Inverse-}\chi^2(\kappa, \lambda).$$

2.2 Kullback–Leibler projection

If p_1 and p_2 are two density functions on \mathbb{R}^d then the *Kullback–Leibler divergence* of p_2 from p_1 is

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

If \mathcal{Q} is a family of univariate density functions then the projection of the univariate density function p onto \mathcal{Q} is

$$\text{proj}_{\mathcal{Q}}[p] \equiv \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(p \| q). \quad (3)$$

A core aspect of expectation propagation is projection of an arbitrary *input* density function onto a particular exponential family. This corresponds to (3) with

$$\mathcal{Q} = \{q(\cdot; \eta) : q(\mathbf{x}; \eta) = \exp\{\mathbf{T}(\mathbf{x})^T \eta - A(\eta)\} h(\mathbf{x}), \quad \eta \in H\}.$$

As explained in Section 2.3 of Kim and Wand (2016), the exponential family Kullback–Leibler problem

$$\eta^* = \underset{\eta \in H}{\text{argmin}} \text{KL}(p \| q(\cdot; \eta))$$

is equivalent to the sufficient statistic moment matching problem

$$\int_{\mathbb{R}^d} \mathbf{T}(\mathbf{x}) \exp\{\mathbf{T}(\mathbf{x})^T \eta^* - A(\eta^*)\} h(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^d} \mathbf{T}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}. \quad (4)$$

Because of (1) we can re-write (4) as

$$(\nabla A)(\eta^*) = \int_{\mathbb{R}^d} \mathbf{T}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$

Then, assuming that the inverse of ∇A is well-defined,

$$\eta^* = (\nabla A)^{-1} \left(\int_{\mathbb{R}^d} \mathbf{T}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \right). \quad (5)$$

Hence, given the \mathbf{T} moments, Kullback–Leibler projection of a density function p onto an exponential family boils down to inversion of ∇A . Section 3 of Wainwright and Jordan (2008) provides a detailed study of exponential families including properties of A and ∇A . An exponential family distribution with the sufficient statistic $\mathbf{T}(\mathbf{x})$ being a $k \times 1$ vector is said to be *regular* if H is an open set in \mathbb{R}^k and *minimal* if there is no $k \times 1$ vector \mathbf{a} and constant $b \in \mathbb{R}$ such that $\mathbf{a}^T \mathbf{T}(\mathbf{x}) = b$ almost surely. Each of the exponential families in Table 1 are regular and minimal. Result 1 provides a summary

of results from Section 3 of Wainwright and Jordan (2008) that is relevant to (5). It depends on:

Definition 1 Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}^k$. Then the set of realizable expectations of f is the set of points $[\tau_1, \dots, \tau_k]^T \in \mathbb{R}^k$ such that there exists a univariate random variable x for which $E\{f(x)\} = [\tau_1, \dots, \tau_k]^T$.

To illustrate the notion of the set of realizable expectations, consider the functions $f_1 : \mathbb{R} \rightarrow \mathbb{R}^2$ and $f_2 : \mathbb{R} \rightarrow \mathbb{R}^2$ given by

$$f_1(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix} \quad \text{and} \quad f_2(x) = \begin{bmatrix} x \\ x^3 \end{bmatrix}.$$

The sets of all realizable expectations of f_1 and f_2 are, respectively

$$\mathfrak{M}_1 \equiv \left\{ \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} : x_2 \geq x_1^2 \right\} \quad \text{and} \quad \mathfrak{M}_2 \equiv \left\{ \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} : \text{sign}(x_1)x_2 \geq |x_1|^3 \right\}.$$

To show that \mathfrak{M}_1 is the set of all realizable expectations of f_1 note that $\mathfrak{M}_1 = \mathfrak{M}_{11} \cup \mathfrak{M}_{12}$ where $\mathfrak{M}_{11} \equiv \{[x_1 \ x_2]^T : x_2 = x_1^2\}$ and $\mathfrak{M}_{12} \equiv \{[x_1 \ x_2]^T : x_2 > x_1^2\}$. Then for any $[x_1 \ x_2] \in \mathfrak{M}_{11}$ we can take x to be the degenerate random variable with probability mass function $p(x) = I(x = x_1)$. For such x , $E\{f_1(x)\} = [x_1 \ x_1^2]^T = [x_1 \ x_2]^T \in \mathfrak{M}_{11}$ which shows that all elements of \mathfrak{M}_{11} are realizable expectations of f_1 . For any $[x_1 \ x_2] \in \mathfrak{M}_{12}$ taking $x \sim N(x_1, x_2 - x_1^2)$ leads to $E\{f_1(x)\} = [x_1 \ x_2]^T$ verifying that all elements of \mathfrak{M}_{12} are realizable by $E\{f_1(x)\}$ for some x . Hence, all entries of \mathfrak{M}_1 are realizable by $E\{f_1(x)\}$ for some x . Values $[x_1 \ x_2]^T \notin \mathfrak{M}_1$ are not realizable because Jensen’s inequality implies that $E(x^2) \geq \{E(x)\}^2$ for any random variable x . Similar arguments can be used to establish that \mathfrak{M}_2 is the set of all realizable expectations of f_2 . Figure 1 shows the sets \mathfrak{M}_1 and \mathfrak{M}_2 .

We are now ready to give the pivotal:

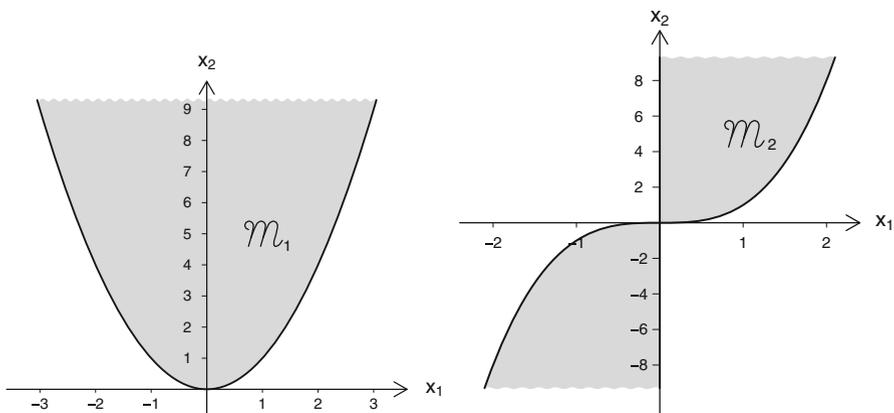


Fig. 1 Left panel: the shaded region is \mathfrak{M}_1 , the set of realizable expectations of f_1 . Right panel: the shaded region is \mathfrak{M}_2 , the set of realizable expectations of f_2

Result 1 (Wainwright and Jordan 2008). Consider a regular and minimal exponential family with k -dimensional sufficient statistic $\mathbf{T}(\mathbf{x})$ and corresponding natural parameter vector $\boldsymbol{\eta}$. Then

- H is a strictly convex subset of \mathbb{R}^k .
- A is a strictly convex and infinitely differentiable function on H .
- ∇A is a one-to-one function.
- The image of ∇A , which we denote by T , is the interior of the set of all realizable expectations of \mathbf{T} .

Result 1 guarantees that $\nabla A : H \rightarrow T$ is a bijective map and that $(\nabla A)^{-1} : T \rightarrow H$ is well-defined.

2.2.1 Normal distribution special case

The Normal distribution is the one of simplest exponential families since ∇A and $(\nabla A)^{-1}$ admit simple closed forms. Firstly, we have

$$\nabla A(\boldsymbol{\eta}) = \begin{bmatrix} -\eta_1/(2\eta_2) \\ (\eta_1^2 - 2\eta_2)/(4\eta_2^2) \end{bmatrix}.$$

It is straightforward to show that the image of H under ∇A is

$$T = \{(\tau_1, \tau_2) : \tau_2 > \tau_1^2\}$$

and the inverse of ∇A is

$$(\nabla A)^{-1}(\boldsymbol{\tau}) = \begin{bmatrix} \tau_1/(\tau_2 - \tau_1^2) \\ -1/\{2(\tau_2 - \tau_1^2)\} \end{bmatrix}.$$

2.2.2 Inverse Chi-Squared distribution special case

For the Inverse Chi-Squared distribution we have

$$\nabla A(\boldsymbol{\eta}) = \begin{bmatrix} \log(-\eta_2) - \text{digamma}(-\eta_1 - 1) \\ (\eta_1 + 1)/\eta_2 \end{bmatrix}$$

where $\text{digamma}(x) \equiv \frac{d}{dx} \log \Gamma(x)$. Determination of the image of H under ∇A is more challenging for the Inverse Chi-Squared distribution. It is aided by Theorem 1 of Kim and Wand (2016) which establishes that $\log -\text{digamma}$ is a bijective map between \mathbb{R}_+ and \mathbb{R}_+ . This leads to

$$T = \{(\tau_1, \tau_2) : \tau_2 > e^{-\tau_1}\}.$$

The inverse of ∇A is

$$(\nabla A)^{-1}(\boldsymbol{\tau}) = \begin{bmatrix} -(\log -\text{digamma})^{-1}(\tau_1 + \log(\tau_2)) - 1 \\ -(\log -\text{digamma})^{-1}(\tau_1 + \log(\tau_2))/\tau_2 \end{bmatrix}.$$

Theorem 1 of Kim and Wand (2016) implies that $(\log -\text{digamma})^{-1}$ is well-defined.

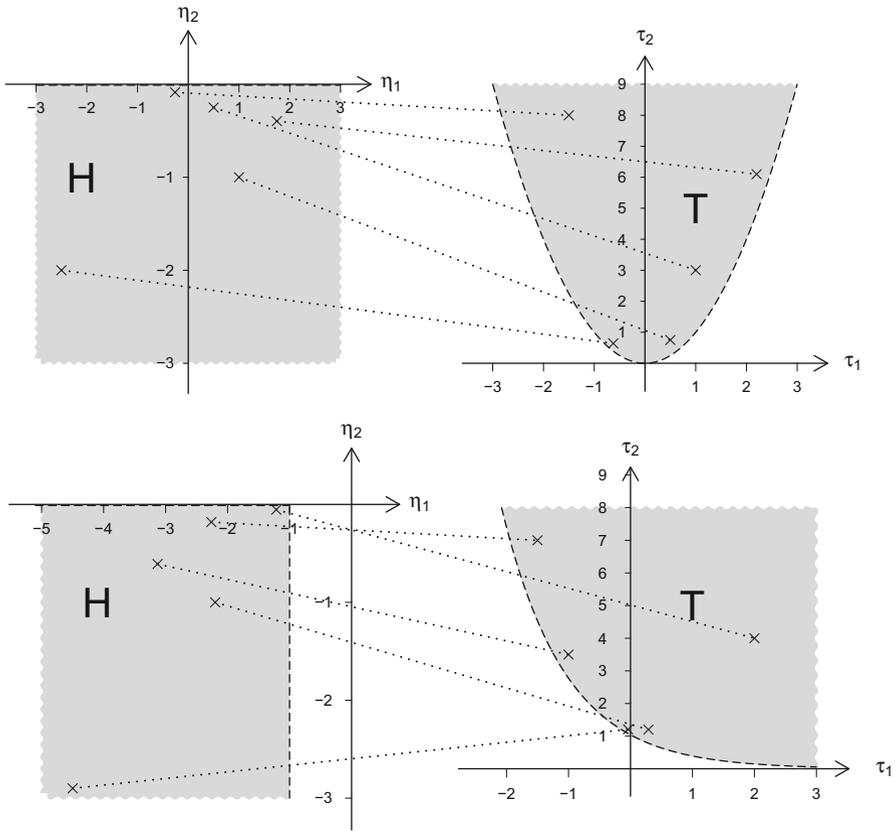


Fig. 2 Upper panel: illustration of the bijective maps between H and T for the Normal exponential family. The crosses and dotted lines depict five example $\eta \in H$ and $\tau = \nabla A(\eta) \in T$ pairs. Since ∇A is a bijective map, the crosses and dotted lines equivalently depict five example $\tau \in T$ and $\eta = (\nabla A)^{-1}(\tau) \in H$ pairs. Lower panel: similar illustration for the Inverse Chi-Squared exponential family

Figure 2 depicts the ∇A and $(\nabla A)^{-1}$ bijective maps between H and T for both the Normal and Inverse Chi-Squared exponential family distributions.

2.3 Factor graphs and factor graph fragments

A *factor graph* is a graphical representation of the factor/argument dependencies of a multivariate function. Even though the concept applies to functions in general, the relevant functions are joint density functions in the context of expectation propagation. As an illustration, consider the Bayesian linear model

$$y|\beta, \sigma^2 \sim N(X\beta, \sigma^2 I),$$

where y is an $n \times 1$ vector of responses, with the following prior distributions on the regression coefficients and error standard deviation:

$$\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}) \text{ and } \sigma \sim \text{Half-Cauchy}(s_{\sigma}),$$

The second prior specification means that σ has prior density function $p(\sigma) = 2/[s_{\sigma}\pi\{1 + (\sigma/s_{\sigma})^2\}]$ for $\sigma > 0$. An equivalent representation of the model, involving the auxiliary variable a , is

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta}, \sigma^2 &\sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\beta}}^2 \mathbf{I}), \\ \sigma^2|a &\sim \text{Inverse} - \chi^2(1, 1/a), \quad a \sim \text{Inverse} - \chi^2(1, 1/s_{\sigma}^2). \end{aligned} \tag{6}$$

We work with this auxiliary variable representation since it aids tractability of expectation propagation. The joint density function of the random variables and random vectors in (6) admits the following factorized form:

$$p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a) = p(\boldsymbol{\beta})p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)p(\sigma^2|a)p(a).$$

Now let \mathbf{x}_i^T be the i th row of \mathbf{X} for $1 \leq i \leq n$. Then a further breakdown of $p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)$ is

$$p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a) = p(\boldsymbol{\beta}) \left\{ \prod_{i=1}^n \int_{-\infty}^{\infty} \delta(\alpha_i - \mathbf{x}_i^T \boldsymbol{\beta}) p(y_i|\alpha_i, \sigma^2) d\alpha_i \right\} p(\sigma^2|a) p(a) \tag{7}$$

where δ is the Dirac delta function and $p(y_i|\alpha_i, \sigma^2) \equiv (2\pi\sigma^2)^{-1/2} \exp\{-(y_i - \alpha_i)^2/(2\sigma^2)\}$. The $\alpha_i \equiv \mathbf{x}_i^T \boldsymbol{\beta}$, $1 \leq i \leq n$, are called *derived variables*. Figure 3 is a factor graph representation of $p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)$ according to the factors that appear in (7). At this point we note that we are not using the conventional factor graph definition here since some of the factors appear inside the integrals in (7). Kim and Wand (2018) introduced the term *derived variable factor graph* to make this distinction. We will simply call it a *factor graph* from now onwards. The circles are called *stochastic nodes* and the rectangles are called *factors*. Both circles and rectangles are *nodes* of the factor graph. We say that a two nodes are *neighbors* of each other if they are joined by an edge.

It is important to note that Fig. 3 is one of many factor graph representations of $p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2, a)$. Other factor graphs arise, for example, via factorizations of $p(\boldsymbol{\beta})$ based on partitions of the $\boldsymbol{\beta}$ vector or taking $p(\sigma^2|a)p(a)$ to be a single factor as a bivariate function of the parameter vector (σ^2, a) . The chosen factor graph matches the type of

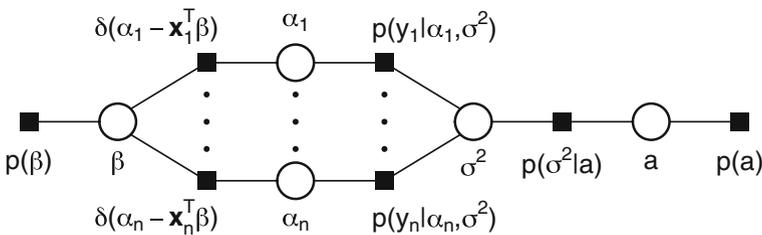


Fig. 3 A factor graph representation of (7)

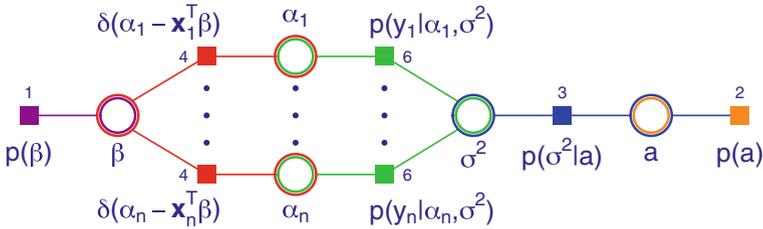


Fig. 4 Fragmentization of the Fig. 3 factor graph. Different colors signify fragments of the same type, and are included in Table 2

mean field restriction imposed on the approximate posterior distribution. For the Fig. 3 factor graph having separate stochastic nodes for β , σ^2 and a corresponds to the following mean field approximation:

$$p(\beta, \sigma^2, a | \mathbf{y}) \approx q(\beta)q(\sigma^2)q(a)$$

for density functions $q(\beta)$, $q(\sigma^2)$ and $q(a)$.

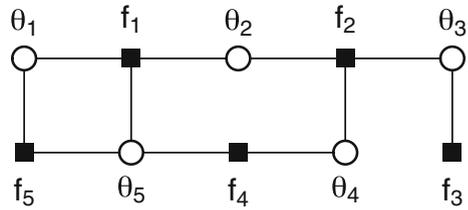
Figure 4 is a representation of Fig. 3 with factor graph *fragments* of the same type identified via color-coding and numbering of the factors. As defined in Wand (2017), a fragment is a sub-graph of a factor graph consisting of a factor and each of its neighboring stochastic nodes.

The five different colors in Fig. 4 correspond to five different fragment types. Some of the fragment types, such as that corresponding to the $p(\beta)$ factor, only appear once in this factor graph. Other types, such as those corresponding to $\delta(\alpha_i - \mathbf{x}_i^T \beta)$, $1 \leq i \leq n$, appear multiple times. Recognition of the recurrence of fragments of the same type in this factor graph and factor graphs for other models is at the core of extension to arbitrarily large models. Wand (2017) demonstrated factor graph fragmentization of variational message passing. Our goal here is to do the same for expectation propagation.

2.4 Expectation propagation

As we discussed in Sect. 1, there are at least two main ways by which expectation propagation is defined for approximate inference in Bayesian models. Throughout this article, we use the message passing on a factor graph version of expectation propagation as developed in Minka (2005) and neatly summarized in Appendix B of Minka and Winn (2008). The preamble of Appendix B of the latter reference contains the sentence “Deterministic factors are not treated specially.” and refers to Dirac delta factors such as the $\delta(\alpha_i - \mathbf{x}_i^T \beta)$ appearing in (7). However, this statement requires adherence to Convention 1, given below. The function $\text{neighbors}(\cdot)$ plays an important role in the algebraic description of the expectation propagation message updates. Consider the illustrative generic form factor graph shown in Fig. 5, corresponding to the joint density function of random vectors $\theta_1, \dots, \theta_5$ according to a particular Bayesian model. Then $\text{neighbours}(1) = \{1, 2, 5\}$ since the factor f_1 is connected by edges to each of θ_1, θ_2 and θ_5 . Similarly, $\text{neighbours}(2) = \{2, 3, 4\}$, $\text{neighbours}(3) = \{3\}$,

Fig. 5 An illustrative generic form factor graph



neighbours(4) = {4, 5} and neighbours(5) = {1, 5}. For general factor graphs with the θ_i and f_j labeling, neighbours(j) is the set of indices of the θ_i that are connected to f_j by an edge.

The message passing version of expectation propagation involves messages passed between neighboring nodes on the factor graph corresponding to the model and mean field restriction. Full details are given in Minka (2005) and Section 3 of Kim and Wand (2016), with the latter reference using the neighbors(\cdot) notation. Based on (54) of Minka (2005), the stochastic node to factor messages are updated according to

$$m_{\theta_i \rightarrow f_j}(\theta_i) \leftarrow \prod_{j' \neq j: i \in \text{neighbours}(j')} m_{f_{j'} \rightarrow \theta_i}(\theta_i) \tag{8}$$

and, based on (83) of Minka (2005), the factor to stochastic node messages updates are

$$m_{f_j \rightarrow \theta_i}(\theta_i) \leftarrow \frac{\text{proj} \left[m_{\theta_i \rightarrow f_j}(\theta_i) \int f_j(\theta_{\text{neighbours}(j)}) \prod_{i' \in \text{neighbours}(j) \setminus \{i\}} m_{\theta_{i'} \rightarrow f_j}(\theta_{i'}) d\theta_{\text{neighbours}(j) \setminus \{i\}} / Z \right]}{m_{\theta_i \rightarrow f_j}(\theta_i)}, \tag{9}$$

where Z is the normalizing factor that ensures that the function of θ_i inside the $\text{proj}[\cdot]$ is a density function. The normalizing factor in (9) involves summation if some of the $\theta_{i'}$ have discrete components. The $\text{proj}[\cdot]$ in (9) denotes Kullback–Leibler projection onto an appropriate exponential family of density functions. However, in Kim and Wand (2016) illustration was done only via a simple example in which all of the stochastic nodes were univariate. In the case of linear models, in which vector parameters are present, some adjustments are necessary to avoid intractable multivariate integrals. The first one is an intrinsically important convention and is now spelt out:

Convention 1. *Derived variable factor graphs are treated as ordinary factor graphs when it comes to applying the message passing expressions (8) and (9).*

In practice, iteration involving (8) and (9) may require some tweaking to achieve convergence. Heskes and Zoeter (2002) provide a careful study of the numerical properties of expectation propagation and show that it can be expressed as a saddle-point problem and argue for the use of *damping* adjustments

$$m_{f_j \rightarrow \theta_i}(\theta_i) \leftarrow m_{f_j \rightarrow \theta_i}(\theta_i)^\epsilon \times \{\text{right-hand side of (9)}\}^{1-\epsilon}, \tag{10}$$

for some $0 \leq \varepsilon < 1$, to enhance convergence. They also propose a more elaborate *double-loop* algorithm with convergence guarantees. Single-loop expectation propagation comes with no convergence guarantees but is considerably easier to implement. In our numerical experiments involving single-loop expectation propagation for semiparametric regression models we have found damping to be both necessary and adequate. Therefore, we build this adjustment into the fragment updates in the next section. The possibility and practicality of double-loop adjustments to fragment updates is an interesting problem for future research.

The full expectation propagation iterative algorithm is:

-
- Initialize all factor to stochastic node messages.
 - Cycle until all factor to stochastic node messages converge:
 - For each factor:
 - Compute the messages passed to the factor using (8).
 - Compute the messages passed from the factor using (9) and (10).
-

Upon convergence the expectation propagation-approximate posterior density function of θ_i is obtained from

$$q^*(\theta_i) \propto \prod_{j:i \in \text{neighbours}(j)} m_{f_j \rightarrow \theta_i}(\theta_i).$$

For the Bayesian linear model example, with the Fig. 3 factor graph, the approximate posterior density functions of β , σ^2 and a are

$$q(\beta) \propto m_{p(\beta) \rightarrow \beta}(\beta) \prod_{i=1}^n m_{\delta(\alpha_i - x_i^T \beta) \rightarrow \beta}(\beta),$$

$$q(\sigma^2) \propto m_{p(\sigma^2|a) \rightarrow \sigma^2}(\sigma^2) \prod_{i=1}^n m_{p(y_i|\alpha_i, \sigma^2) \rightarrow \sigma^2}(\sigma^2) \quad \text{and}$$

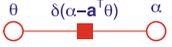
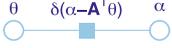
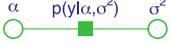
$$q(a) \propto m_{p(\sigma^2|a) \rightarrow a}(a) m_{p(a) \rightarrow a}(a).$$

3 Fragmentization for generalized, linear and mixed models

In principle, factor graph fragmentation of expectation propagation applies to any graphical model. In this section we provide illustration of fragmentation for generalized, linear and mixed models. This is rich class of models that also includes may semiparametric regression models, such as generalized additive models, via the use of mixed model representations of penalized splines (e.g. Ruppert et al. 2009). By the end of this section we will have derived and catalogued several important fragments for expectation propagation-based approximate Bayesian inference and paved the way for similar work for other model families.

Each of the generalized, linear and mixed models dealt with in Kim and Wand (2018) can be handled with nine distinct fragment types, which are listed in Table 2. The

Table 2 Fundamental factor graph fragments for expectation propagation fitting of generalized, linear and mixed models

Fragment name	Diagram	Distributional statement
1. Gaussian prior		$\theta \sim N(\mu_\theta, \Sigma_\theta)$
2. Inverse Wishart prior		$\Theta \sim \text{Inverse-Wishart}(\kappa_\Theta, \Lambda_\Theta)$
3. Iterated Inverse Chi-Squared		$\sigma^2 a \sim \text{Inverse-}\chi^2(v, 1/a)$
4. Linear combination derived variable		$\alpha \equiv a^T \theta$
5. Multivariate linear combination derived variable		$\alpha \equiv A^T \theta$
6. Gaussian		$y \alpha, \sigma^2 \sim N(\alpha, \sigma^2)$
7. Logistic likelihood		$y \alpha \sim \text{Bernoulli}(\text{logit}^{-1}(\alpha))$
8. Probit likelihood		$y \alpha \sim \text{Bernoulli}(\Phi(\alpha))$
9. Poisson likelihood		$y \alpha \sim \text{Poisson}(\exp(\alpha))$

message updates for each fragment type only needs to be derived once. Each subsection deals with the required derivation and summarizes the updates as an algorithm. For a software suite that uses expectation propagation to fit generalized, linear and mixed models the fragment only needs to be implemented once. We now work through each of the Table 2 fragments in turn.

The algorithms use the matrix functions vec and its inverse vec^{-1} which we define here. If A is $d \times d$ matrix then $\text{vec}(A)$ is the $d^2 \times 1$ vector obtained by stacking the columns of A underneath each other in order from left to right. if a is a $d^2 \times 1$ vector then $\text{vec}^{-1}(a)$ is the $d \times d$ matrix formed from listing the entries of a in a column-wise fashion in order from left to right and is the usual function inverse when the domain of vec is restricted to square matrices.

The following shorthand is used throughout this section:

$$a \xleftarrow{\varepsilon} b \text{ denotes } a \leftarrow \varepsilon a + (1 - \varepsilon) b.$$

3.1 Gaussian prior fragment

The Gaussian prior fragment arises from the following prior distribution specification:

$$\theta \sim N(\mu_\theta, \Sigma_\theta)$$

for user-specified hyperparameters μ_θ and Σ_θ . The fragment factor is

$$p(\theta) = (2\pi)^{-d_\theta/1} |\Sigma_\theta|^{-1/2} \exp \left\{ -\frac{1}{2}(\theta - \mu_\theta)^T \Sigma_\theta^{-1} (\theta - \mu_\theta) \right\}.$$

We assume that

$$\text{all messages passed to } \theta \text{ from factors outside of the fragment are in the Multivariate Normal family.} \tag{11}$$

The message from $p(\theta)$ to θ takes the form

$$m_{p(\theta) \rightarrow \theta}(\theta) = \exp \left\{ \left[\begin{array}{c} \theta \\ \text{vec}(\theta\theta^T) \end{array} \right]^T \eta_{p(\theta) \rightarrow \theta} \right\}.$$

Algorithm 1 provides the natural parameter update for this simple fragment. The derivation of Algorithm 1 is given in Section S.2.1 of the online supplement.

Algorithm 1 *The input, update and output of the Gaussian prior fragment.*

Hyperparameter Inputs: $\mu_\theta, \Sigma_\theta$.

Update:

$$\eta_{p(\theta) \rightarrow \theta} \leftarrow \left[\begin{array}{c} \Sigma_\theta^{-1} \mu_\theta \\ -\frac{1}{2} \text{vec}(\Sigma_\theta^{-1}) \end{array} \right]$$

Parameter Output: $\eta_{p(\theta) \rightarrow \theta}$.

3.2 Inverse Wishart prior fragment

Let Θ be a $d^\Theta \times d^\Theta$ symmetric positive definite random matrix. The prior specification

$$\Theta \sim \text{Inverse-Wishart}(\kappa_\Theta, \Lambda_\Theta)$$

leads to a factor graph fragment with factor

$$p(\Theta) = \mathcal{C}_{d^\Theta, \kappa_\Theta}^{-1} |\Lambda_\Theta|^{\kappa_\Theta/2} |\Theta|^{-(\kappa_\Theta + d^\Theta + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_\Theta \Theta^{-1}) \right\} I(\Theta \text{ symmetric and positive definite}).$$

where $\mathcal{C}_{d^\Theta, \kappa_\Theta}$ is defined via (2). The message from $p(\Theta)$ to Θ takes the form

$$m_{p(\Theta) \rightarrow \Theta}(\Theta) = \exp \left\{ \left[\begin{array}{c} \log |\Theta| \\ \text{vec}(\Theta^{-1}) \end{array} \right]^T \eta_{p(\Theta) \rightarrow \Theta} \right\}.$$

Algorithm 2 gives the $\eta_{p(\Theta) \rightarrow \Theta}$ update based on hyperparameter inputs κ_Θ and Λ_Θ .

Algorithm 2 *The input, update and output of the Inverse Wishart prior fragment.*

Hyperparameter Inputs: κ_{Θ} , Λ_{Θ} .

Update:

$$\eta_{p(\Theta) \rightarrow \Theta} \leftarrow \begin{bmatrix} -\frac{1}{2}(\kappa_{\Theta} + d^{\Theta} + 1) \\ -\frac{1}{2}\text{vec}(\Lambda_{\Theta}) \end{bmatrix}$$

Parameter Output: $\eta_{p(\Theta) \rightarrow \Theta}$.

A derivation of Algorithm 2 is given in Section S.2.2 of the online supplement.

3.3 Iterated Inverse Chi-Squared fragment

This fragment arises from the following distributional fact (e.g. Wand et al. (2011), Result 5):

$$\begin{aligned} \sigma^2 | a \sim \text{Inverse-}\chi^2(\nu, \nu/a) \quad \text{and} \quad a \sim \text{Inverse-}\chi^2(1, 1/A^2) \\ \text{implies} \quad \sigma \sim \text{Half-}t(A, \nu) \end{aligned} \quad (12)$$

where $x \sim \text{Half-}t(A, \nu)$ if and only if

$$p(x) = \frac{2\Gamma\left(\frac{\nu+1}{2}\right) I(x > 0)}{\sqrt{\pi\nu} \Gamma(\nu/2) A \{1 + (x/A)^2/\nu\}^{(\nu+1)/2}}.$$

The advantage of fact (12) is that non-informative priors within the Half- t family can be imposed on standard deviation parameters using messages within the Inverse Chi-Squared family.

The fragment factor is

$$p(\sigma^2 | a) = \frac{\{v/(2a)\}^{v/2}}{\Gamma(v/2)} (\sigma^2)^{-(v/2)-1} \exp\{-v/(2a\sigma^2)\} I(\sigma^2 > 0) I(a > 0)$$

and it is assumed that:

all messages passed to σ^2 from factors outside of
the fragment are in the Inverse Chi-Squared family and
all messages passed to a from factors outside of
the fragment are also in the Inverse Chi-Squared family. (13)

The messages from the factor to its neighboring stochastic nodes are

$$m_{p(\sigma^2 | a) \rightarrow \sigma^2}(\sigma^2) = \exp \left\{ \begin{bmatrix} \log(\sigma^2) \\ 1/\sigma^2 \end{bmatrix}^T \eta_{p(\sigma^2 | a) \rightarrow \sigma^2} \right\}$$

and

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

Algorithm 3 provides the updates of the natural parameters of these messages given messages from outside the fragment. The function G^{IG3} is defined in Section S.1.3.

Algorithm 3 *The inputs, updates and outputs of the iterated Inverse Chi-Squared fragment.*

Data Input: $\nu > 0, 0 \leq \varepsilon < 1.$

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

Updates:

$$\begin{aligned} \boldsymbol{\eta}_{p(\sigma^2|a)\rightarrow \sigma^2} &\stackrel{\varepsilon}{\leftarrow} G^{IG3} \left(\boldsymbol{\eta}_{\sigma^2\rightarrow p(\sigma^2|a)}, \boldsymbol{\eta}_{a\rightarrow p(\sigma^2|a)}; \nu + 2, \nu \right) \\ \boldsymbol{\eta}_{p(\sigma^2|a)\rightarrow a} &\stackrel{\varepsilon}{\leftarrow} G^{IG3} \left(\boldsymbol{\eta}_{a\rightarrow p(\sigma^2|a)}, \boldsymbol{\eta}_{\sigma^2\rightarrow p(\sigma^2|a)}; \nu, \nu \right) \end{aligned}$$

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

3.4 Linear combination derived variable fragment

The linear combination derived variable fragment corresponds to equating a scalar variable α with a linear combination $\mathbf{a}^T \boldsymbol{\theta}$, where $\boldsymbol{\theta}$ is a Multivariate Normal random vector. If g is a general function that depends on the linear combination form $\mathbf{a}^T \boldsymbol{\theta}$ and other variables, denoted by \mathbf{o} , then the derived variable α arises from the equality:

$$g(\mathbf{a}^T \boldsymbol{\theta}; \mathbf{o}) = \int_{-\infty}^{\infty} \delta(\alpha - \mathbf{a}^T \boldsymbol{\theta}) g(\alpha; \mathbf{o}) d\alpha. \tag{14}$$

where δ is the Dirac delta function. Under Convention 1 given in Sect. 2.4, the integral sign is ignored when it comes to applying the expectation propagation updates (8) and (9). We assume that:

$$\begin{aligned} &\text{all messages passed to } \alpha \text{ from factors outside of} \\ &\text{the fragment are in the Univariate Normal family} \\ &\text{and all messages passed to } \boldsymbol{\theta} \text{ from factors outside} \\ &\text{of the fragment are in the Multivariate Normal family.} \end{aligned} \tag{15}$$

The function $\delta(\alpha - \mathbf{a}^T \boldsymbol{\theta})$ is the factor for this fragment. According to conjugacy restrictions, messages passed from $\delta(\alpha - \mathbf{a}^T \boldsymbol{\theta})$ to α and $\boldsymbol{\theta}$ take the forms

$$m_{\delta(\alpha - a^T \theta) \rightarrow \alpha}(\alpha) = \exp \left\{ \begin{bmatrix} \alpha \\ \alpha^2 \end{bmatrix}^T \eta_{\delta(\alpha - a^T \theta) \rightarrow \alpha} \right\} \quad (16)$$

and

$$m_{\delta(\alpha - a^T \theta) \rightarrow \theta}(\theta) = \exp \left\{ \begin{bmatrix} \theta \\ \text{vec}(\theta \theta^T) \end{bmatrix}^T \eta_{\delta(\alpha - a^T \theta) \rightarrow \theta} \right\}.$$

Algorithm 4 provides the updates to the natural parameter vectors

$$\eta_{\delta(\alpha - a^T \theta) \rightarrow \alpha} \quad \text{and} \quad \eta_{\delta(\alpha - a^T \theta) \rightarrow \theta}$$

given inputs

$$\eta_{\alpha \rightarrow \delta(\alpha - a^T \theta)} \quad \text{and} \quad \eta_{\theta \rightarrow \delta(\alpha - a^T \theta)}.$$

It uses the notation

$$\begin{aligned} (\eta_{\theta \rightarrow \delta(\alpha - a^T \theta)})_1 &\equiv \text{vector containing the first } d^\theta \text{ entries of } \eta_{\theta \rightarrow \delta(\alpha - a^T \theta)} \\ \text{and } (\eta_{\theta \rightarrow \delta(\alpha - a^T \theta)})_2 &\equiv \text{vector containing the remaining } (d^\theta)^2 \text{ entries of } \eta_{\theta \rightarrow \delta(\alpha - a^T \theta)} \end{aligned}$$

where d^θ is the number of entries in θ . The derivations of these updates are given in Section S.2.5 of the online supplement. The updates have analogies with those presented in Herbrich (2005).

Algorithm 4 *The inputs, updates and outputs of the linear combination derived variable fragment.*

Data Input: a (vector having the same dimension as θ), $0 \leq \varepsilon < 1$.

Parameter Inputs: $\eta_{\delta(\alpha - a^T \theta) \rightarrow \alpha}$, $\eta_{\delta(\alpha - a^T \theta) \rightarrow \theta}$, $\eta_{\alpha \rightarrow \delta(\alpha - a^T \theta)}$, $\eta_{\theta \rightarrow \delta(\alpha - a^T \theta)}$.

Updates:

$$\begin{aligned} \omega &\leftarrow \left\{ \text{vec}^{-1} \left((\eta_{\theta \rightarrow \delta(\alpha - a^T \theta)})_2 \right) \right\}^{-1} a \\ \eta_{\delta(\alpha - a^T \theta) \rightarrow \alpha} &\leftarrow^{\varepsilon} \frac{1}{\omega^T a} \begin{bmatrix} \omega^T (\eta_{\theta \rightarrow \delta(\alpha - a^T \theta)})_1 \\ 1 \end{bmatrix} \\ \eta_{\delta(\alpha - a^T \theta) \rightarrow \theta} &\leftarrow^{\varepsilon} \begin{bmatrix} a (\eta_{\alpha \rightarrow \delta(\alpha - a^T \theta)})_1 \\ \text{vec}(a a^T) (\eta_{\alpha \rightarrow \delta(\alpha - a^T \theta)})_2 \end{bmatrix} \end{aligned}$$

Parameter Outputs: $\eta_{\delta(\alpha - a^T \theta) \rightarrow \alpha}$, $\eta_{\delta(\alpha - a^T \theta) \rightarrow \theta}$.

3.5 Multivariate linear combination derived variable fragment

Now consider the following bivariate extension of (14):

$$g(\mathbf{a}_1^T \boldsymbol{\theta}, \mathbf{a}_2^T \boldsymbol{\theta}; \boldsymbol{o}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta(\alpha_1 - \mathbf{a}_1^T \boldsymbol{\theta}) \delta(\alpha_2 - \mathbf{a}_2^T \boldsymbol{\theta}) g(\alpha_1, \alpha_2; \boldsymbol{o}) d\alpha_1 d\alpha_2, \tag{17}$$

where the primary argument of the function g is now bivariate. The established result for the Dirac delta function applied to bivariate arguments leads to the equivalent form for the right-hand side of (17) taking the form:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta \left(\begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} - \mathbf{A}^T \boldsymbol{\theta} \right) g(\alpha_1, \alpha_2; \boldsymbol{o}) d\alpha_1 d\alpha_2 \quad \text{where } \mathbf{A} \equiv [\mathbf{a}_1 \ \mathbf{a}_2].$$

It follows that

$$\boldsymbol{\alpha} \equiv \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix}$$

is a bivariate derived variable corresponding to the multivariate linear combination $\mathbf{A}^T \boldsymbol{\theta}$.

In the most general case, $\boldsymbol{\theta}$ and $\boldsymbol{\alpha}$ are, respectively, $d^\theta \times 1$ and $d^\alpha \times 1$ vectors and \mathbf{A} is a $d^\theta \times d^\alpha$ matrix. The fragment factor is $\delta(\boldsymbol{\alpha} - \mathbf{A}^T \boldsymbol{\theta})$ and the message given in (16) generalizes to

$$\eta_{\delta(\boldsymbol{\alpha} - \mathbf{A}^T \boldsymbol{\theta}) \rightarrow \boldsymbol{\alpha}}(\boldsymbol{\alpha}) = \exp \left\{ \left[\begin{bmatrix} \boldsymbol{\alpha} \\ \text{vec}(\boldsymbol{\alpha} \boldsymbol{\alpha}^T) \end{bmatrix} \right]^T \eta_{\delta(\boldsymbol{\alpha} - \mathbf{A}^T \boldsymbol{\theta}) \rightarrow \boldsymbol{\alpha}} \right\}.$$

Algorithm 5 lists the natural parameter updates. Their derivations are given in Section S.2.5 of the online supplement.

Algorithm 5 *The inputs, updates and outputs of the multivariate linear combination derived variable fragment.*

Data Input: A (matrix with number of columns matching the dimension of θ), $0 \leq \varepsilon < 1$.

Parameter Inputs: $\eta_{\delta(\alpha-A^T\theta)\rightarrow\alpha}, \eta_{\delta(\alpha-A^T\theta)\rightarrow\theta}, \eta_{\alpha\rightarrow\delta(\alpha-A^T\theta)}, \eta_{\theta\rightarrow\delta(\alpha-A^T\theta)}$.

Updates:

$$\begin{aligned} \Omega &\leftarrow \left\{ \text{vec}^{-1} \left(\left(\eta_{\theta\rightarrow\delta(\alpha-A^T\theta)} \right)_2 \right) \right\}^{-1} A \\ \eta_{\delta(\alpha-A^T\theta)\rightarrow\alpha} &\leftarrow^{\varepsilon} \begin{bmatrix} (\Omega^T A)^{-1} \Omega^T \left(\eta_{\theta\rightarrow\delta(\alpha-A^T\theta)} \right)_1 \\ \text{vec} \left((\Omega^T A)^{-1} \right) \end{bmatrix} \\ \eta_{\delta(\alpha-A^T\theta)\rightarrow\theta} &\leftarrow^{\varepsilon} \begin{bmatrix} A \left(\eta_{\alpha\rightarrow\delta(\alpha-A^T\theta)} \right)_1 \\ (A \otimes A) \left(\eta_{\alpha\rightarrow\delta(\alpha-A^T\theta)} \right)_2 \end{bmatrix} \end{aligned}$$

Parameter Outputs: $\eta_{\delta(\alpha-A^T\theta)\rightarrow\alpha}, \eta_{\delta(\alpha-A^T\theta)\rightarrow\theta}$.

Note that Algorithm 5 is a generalization of Algorithm 4. Therefore, from a strict mathematical standpoint, Algorithm 4 is unnecessary. However, since ordinary linear combinations are common in expectation propagation fitting of linear models we feel that it is worth having a separate fragment and algorithm for this special case.

3.6 Gaussian fragment

The Gaussian fragment corresponds to the specification

$$y|\alpha, \sigma^2 \sim N(\alpha, \sigma^2).$$

The fragment's factor is

$$p(y|\alpha, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp\{-(y - \alpha)^2/(2\sigma^2)\}$$

which, as a function of α , is in the Normal family and, as a function of σ^2 , is in the Inverse Chi-Squared family. Exponential family constraint considerations then lead to the following assumption for the Gaussian fragment:

all messages passed to α from factors outside of
the fragment are in the Univariate Normal family
and all messages passed to σ^2 from factors outside
of the fragment are in the Inverse Chi-Squared family. (18)

The messages from $p(y|\alpha, \sigma^2)$ take the forms

$$m_{p(y|\alpha, \sigma^2)\rightarrow\alpha}(\alpha) = \exp \left\{ \begin{bmatrix} \alpha \\ \alpha^2 \end{bmatrix}^T \eta_{p(y|\alpha, \sigma^2)\rightarrow\alpha} \right\}$$

and

$$m_{p(y|\alpha, \sigma^2) \rightarrow \sigma^2}(\sigma^2) = \exp \left\{ \begin{bmatrix} \log(\sigma^2) \\ 1/\sigma^2 \end{bmatrix}^T \boldsymbol{\eta}_{p(y|\alpha, \sigma^2) \rightarrow \sigma^2} \right\}$$

with natural parameters updated according to Algorithm 6. The functions G^N and G^{IG3} are defined in Section S.1.3. Algorithm 6’s derivation is given in Section S.2.6.

Algorithm 6 *The inputs, updates and outputs of the Gaussian fragment.*

Data Input: $y \in \mathbb{R}$, $0 \leq \varepsilon < 1$.

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

Update:

$$\begin{aligned} \boldsymbol{\eta}_{p(y|\alpha, \sigma^2) \rightarrow \alpha} &\stackrel{\varepsilon}{\leftarrow} G^N \left(\boldsymbol{\eta}_{\alpha \rightarrow p(y|\alpha, \sigma^2)}, \boldsymbol{\eta}_{\sigma^2 \rightarrow p(y|\alpha, \sigma^2)}; \begin{bmatrix} 1 \\ y \\ y^2 \end{bmatrix} \right) \\ \boldsymbol{\eta}_{p(y|\alpha, \sigma^2) \rightarrow \sigma^2} &\stackrel{\varepsilon}{\leftarrow} G^{IG1} \left(\boldsymbol{\eta}_{\sigma^2 \rightarrow p(y|\alpha, \sigma^2)}, \boldsymbol{\eta}_{\alpha \rightarrow p(y|\alpha, \sigma^2)}; \begin{bmatrix} 1 \\ y \\ y^2 \end{bmatrix} \right) \end{aligned}$$

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

3.7 Logistic likelihood fragment

The logistic likelihood fragment corresponds to the specification

$$y|\alpha \sim \text{Bernoulli}\{\text{logit}^{-1}(\alpha)\}.$$

The factor of the fragment is

$$p(y|\alpha) = \exp\{y\alpha - \log(1 + e^\alpha)\}.$$

We assume that:

$$\begin{aligned} &\text{all messages passed to } \alpha \text{ from other factors are} \\ &\text{within the Univariate Normal exponential family.} \end{aligned} \tag{19}$$

Conjugacy then dictates that

$$m_{p(y|\alpha) \rightarrow \alpha}(\alpha) = \exp \left\{ \begin{bmatrix} \alpha \\ \alpha^2 \end{bmatrix}^T \boldsymbol{\eta}_{p(y|\alpha) \rightarrow \alpha} \right\}. \tag{20}$$

Algorithm 7 provides the update to the natural parameter vector

$$\eta_{p(y|\alpha)\rightarrow\alpha} \text{ based on input } \eta_{\alpha\rightarrow p(y|\alpha)}$$

and depends on the function H_{logistic} defined at equation (S.3) in the online supplement. Its derivation is given in Section S.2.7 of the online supplement.

3.8 Probit likelihood fragment

The probit likelihood fragment corresponds to the specification

$$y|\alpha \sim \text{Bernoulli}(\Phi(\alpha)).$$

Algorithm 7 *The inputs, updates and outputs of the logistic likelihood fragment.*

Data Input: $y \in \{0, 1\}$, $0 \leq \varepsilon < 1$.

Parameter Inputs: $\eta_{p(y|\alpha)\rightarrow\alpha}$, $\eta_{\alpha\rightarrow p(y|\alpha)}$.

Update:

$$\eta_{p(y|\alpha)\rightarrow\alpha} \leftarrow^{\varepsilon} H_{\text{logistic}}(\eta_{\alpha\rightarrow p(y|\alpha)}; y)$$

Parameter Output: $\eta_{p(y|\alpha)\rightarrow\alpha}$.

The factor of the fragment is

$$p(y|\alpha) = \exp [y \log\{\Phi(\alpha)\} + (1 - y) \log\{1 - \Phi(\alpha)\}].$$

As for the logistic likelihood fragment, we also assume (19) which implies that $m_{p(y|\alpha)\rightarrow\alpha}(\alpha)$ also takes the form (20). The fragment update is given in Algorithm 8, with justification deferred to Section S.2.8 of the online supplement. The function H_{probit} is defined in Section S.1.3 of the online supplement. Note that H_{probit} has the advantage of admitting a closed form expression. This is not the case for H_{logistic} and numerical integration is required for its evaluation.

Algorithm 8 *The inputs, updates and outputs of the probit likelihood fragment.*

Data Input: $y \in \{0, 1\}$, $0 \leq \varepsilon < 1$.

Parameter Inputs: $\eta_{p(y|\alpha) \rightarrow \alpha}$, $\eta_{\alpha \rightarrow p(y|\alpha)}$.

Update:

$$\eta_{p(y|\alpha) \rightarrow \alpha} \stackrel{\varepsilon}{\leftarrow} H_{\text{probit}}(\eta_{\alpha \rightarrow p(y|\alpha)}; y)$$

Parameter Outputs: $\eta_{p(y|\alpha) \rightarrow \alpha}$.

3.9 Poisson likelihood fragment

The Poisson likelihood fragment matches

$$y|\alpha \sim \text{Poisson}\{\exp(\alpha)\}$$

and the factor of the fragment is

$$p(y|\alpha) = \exp\{y\alpha - e^\alpha - \log(y!)\}.$$

As for the logistic and Poisson likelihood fragments, we also assume (19) which implies that $m_{p(y|\alpha) \rightarrow \alpha}(\alpha)$ also takes the form (20). The fragment update is given in Algorithm 9 with the H_{Poisson} function defined at equation (S.3) of the online supplement.

Section S.2.9 of the online supplement contains justification of Algorithm 9.

Algorithm 9 *The inputs, updates and outputs of the Poisson likelihood fragment.*

Data Input: $y \in \{0, 1, 2, \dots\}$, $0 \leq \varepsilon < 1$.

Parameter Inputs: $\eta_{p(y|\alpha) \rightarrow \alpha}$, $\eta_{\alpha \rightarrow p(y|\alpha)}$.

Update:

$$\eta_{p(y|\alpha) \rightarrow \alpha} \stackrel{\varepsilon}{\leftarrow} H_{\text{Poisson}}(\eta_{\alpha \rightarrow p(y|\alpha)}; y)$$

Parameter Outputs: $\eta_{p(y|\alpha) \rightarrow \alpha}$.

4 Illustration

We now provide illustration via a generalized additive mixed model analysis. The data are from the Indonesian Children’s Health Study (Sommer 1982), corresponding to a cohort of 275 Indonesian children who are repeatedly examined. The response variable is

$$y_{ij} = \begin{cases} 1, & \text{if respiratory infection present in the } i\text{th child at the } j\text{th examination,} \\ 0, & \text{otherwise} \end{cases} \tag{21}$$

for $1 \leq i \leq m$ and $1 \leq j \leq n_i$. For these data note that $m = 275$ and the $n_i \in \{1, \dots, 6\}$. Potential predictor variables are age, indicator of vitamin A deficiency, indicator of being female, height, indicator of being stunted and indicators for the number of clinic visits for each child. We let a_{ij} denote the age in years of the i th child at the j th examination. Consider the following Bayesian generalized additive mixed model:

$$\begin{aligned}
 y_{ij} | \beta_0, \boldsymbol{\beta}_x, \beta_{\text{spl}}, \mathbf{u}_{\text{grp}}, \mathbf{u}_{\text{spl}} &\stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left(\text{logit}^{-1} (\beta_0 + u_{\text{grp},i} + \boldsymbol{\beta}_x^T \mathbf{x}_{ij} + f(a_{ij})) \right), \\
 f(a_{ij}) &\equiv \beta_{\text{spl}} a_{ij} + \sum_{k=1}^K u_{\text{spl},k} z_k(a_{ij}) \text{ is a low-rank smoothing spline in } a_{ij}, \\
 &\text{where } \{z_k(\cdot) : 1 \leq k \leq K\} \text{ is a suitable spline basis,} \\
 \boldsymbol{\beta} \equiv \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta}_x \\ \beta_{\text{spl}} \end{bmatrix} &\sim N(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta), \quad \mathbf{u} \equiv \begin{bmatrix} \mathbf{u}_{\text{grp}} \\ \mathbf{u}_{\text{spl}} \end{bmatrix} \Big| \sigma_{\text{grp}}^2, \sigma_{\text{spl}}^2 \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \sigma_{\text{grp}}^2 \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \sigma_{\text{spl}}^2 \mathbf{I}_K \end{bmatrix} \right), \\
 \sigma_{\text{grp}}^2 | a_{\text{grp}} &\sim \text{Inverse-}\chi^2(1, 1/a_{\text{grp}}), \quad \sigma_{\text{spl}}^2 | a_{\text{spl}} \sim \text{Inverse-}\chi^2(1, 1/a_{\text{spl}}), \\
 a_{\text{grp}} &\sim \text{Inverse-}\chi^2(1, 1/s_{\text{grp}}^2), \quad a_{\text{spl}} \sim \text{Inverse-}\chi^2(1, 1/s_{\text{spl}}^2).
 \end{aligned} \tag{22}$$

The ‘grp’ and ‘spl’ subscripting indicates whether the random effect vector and corresponding variance parameter is for the random subject intercept or for spline coefficients in the non-linear function of age. Let \mathbf{y} denote the $N \times 1$ vector containing the y_{ij} , where $N \equiv \sum_{i=1}^m n_i$. Despite the common use of double subscript notation as in (21), it is more convenient to label the entries of \mathbf{y} with a single subscript when it comes to fitting via expectation propagation. To avoid a notational clash we use y_ℓ^s , $1 \leq \ell \leq N$, to denote the ℓ th entry of \mathbf{y} . Let d^β be the number of rows in $\boldsymbol{\beta}$. For the Indonesian Children’s Health Study Data application $d^\beta = 11$. Then let \mathbf{X} by the $N \times d^\beta$ matrix containing the predictor data. The random effects design matrix is $\mathbf{Z} = [\mathbf{Z}_{\text{grp}} \ \mathbf{Z}_{\text{spl}}]$ where

$$\mathbf{Z}_{\text{grp}} \equiv \text{blockdiag}(\mathbf{1}_{n_i}) \quad \text{and} \quad \mathbf{Z}_{\text{spl}} \equiv \begin{bmatrix} z_k(a_{ij}) \\ 1 \leq k \leq K \end{bmatrix}_{1 \leq j \leq n_i, 1 \leq i \leq m}.$$

Then the likelihood can be written as

$$y_\ell^s | \boldsymbol{\beta}, \mathbf{u} \stackrel{\text{ind.}}{\sim} \text{Bernoulli} \left(\text{logit}^{-1} ((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_\ell) \right), \quad 1 \leq \ell \leq N.$$

Next, let $C \equiv [X \ Z]$ so that

$$X\beta + Zu = C \begin{bmatrix} \beta \\ u \end{bmatrix}$$

and let c_ℓ^T be the ℓ th row of C . Let e_r be the $(m + K) \times 1$ vector with r th entry equal to 1 and zeroes elsewhere for $1 \leq r \leq m + K$. Lastly, let E_{d^β} be the $(d^\beta + m + K) \times d^\beta$ matrix with the $d^\beta \times d^\beta$ identity matrix at the top and all other entries equal to zero. The joint density function of all random variables in the model is

$$\begin{aligned} p(y, \beta, u, \sigma_{\text{grp}}^2, \sigma_{\text{spl}}^2, a_{\text{grp}}, a_{\text{spl}}) &= p(y|\beta, u) p(\beta) p(u|\sigma_{\text{grp}}^2, \sigma_{\text{spl}}^2) p(\sigma_{\text{grp}}^2|a_{\text{grp}}) p(\sigma_{\text{spl}}^2|a_{\text{spl}}) p(a_{\text{grp}}) p(a_{\text{spl}}) \\ &= p(\beta) \left\{ \prod_{\ell=1}^N p(y_\ell^s | \beta, u) \right\} \left\{ \prod_{i=1}^m p(u_{\text{grp},i} | \sigma_{\text{grp}}^2) \right\} \left\{ \prod_{k=1}^K p(u_{\text{spl},k} | \sigma_{\text{spl}}^2) \right\} p(\sigma_{\text{grp}}^2|a_{\text{grp}}) \\ &\quad \times p(\sigma_{\text{spl}}^2|a_{\text{spl}}) p(a_{\text{grp}}) p(a_{\text{spl}}) \\ &= \left\{ \int_{\mathbb{R}^{d^\beta}} p(\tilde{\beta}) \delta \left(\tilde{\beta} - E_{d^\beta}^T \begin{bmatrix} \beta \\ u \end{bmatrix} \right) d\tilde{\beta} \right\} \left\{ \prod_{\ell=1}^N \int_{-\infty}^{\infty} p(y_\ell^s | \alpha_\ell) \delta \left(\alpha_\ell - c_\ell^T \begin{bmatrix} \beta \\ u \end{bmatrix} \right) d\alpha_\ell \right\} \\ &\quad \times \left\{ \prod_{i=1}^m \int_{-\infty}^{\infty} p(\tilde{u}_{\text{grp},i} | \sigma_{\text{grp}}^2) \delta \left(\tilde{u}_{\text{grp},i} - e_{d^\beta+i}^T \begin{bmatrix} \beta \\ u \end{bmatrix} \right) d\tilde{u}_{\text{grp},i} \right\} p(\sigma_{\text{grp}}^2|a_{\text{grp}}) p(a_{\text{grp}}) \\ &\quad \times \left\{ \prod_{k=1}^K \int_{-\infty}^{\infty} p(\tilde{u}_{\text{spl},k} | \sigma_{\text{spl}}^2) \delta \left(\tilde{u}_{\text{spl},k} - e_{d^\beta+m+k}^T \begin{bmatrix} \beta \\ u \end{bmatrix} \right) d\tilde{u}_{\text{spl},k} \right\} p(\sigma_{\text{spl}}^2|a_{\text{spl}}) p(a_{\text{spl}}). \end{aligned} \tag{23}$$

Figure 6 is the derived variable factor graph corresponding to the representation of the joint density function given in (23). All of the fragments in Fig. 6 are versions of fundamental fragments listed in Table 2 and are color-coded and numbered accordingly. Expectation propagation inference for this model and data involves iteratively passing messages between neighboring nodes on the Fig. 6 factor graph. The parameter updates for the factor to stochastic node messages are given by the relevant algorithms in Sect. 3. The stochastic node to factor message parameter updates are a simple consequence of (8).

We fit (22) using 1,000 iterations of expectation propagation message passing on the factor graph of Fig. 6. We also conducted Markov chain Monte Carlo fitting via the function `stan()` in the R package `rstan` (Guo et al. 2017), which interfaces the Stan language (Carpenter et al. 2017), with a warmup size of 50,000 and a retained sample size of 1,000,000. The hyperparameters were set to $\mu_\beta = \mathbf{0}$, $\Sigma_\beta = 10^{10}I$, $\sigma_{\text{grp}} = \sigma_{\text{spl}} = 10^5$ with continuous variables standardized for the analyses and then results transformed to correspond to the original units. Figure 7 compares the Bayesian inference arising from the two approaches. The first three rows compare the expectation and Markov chain Monte Carlo approximate posterior density functions for the fixed effects parameters. The last row contains similar comparisons for the variance parameters and the low-rank smoothing spline fits for the non-linear age effect. The estimated probability functions are such that all other predictors are set at their average values, and are accompanied by pointwise 95% credible intervals.

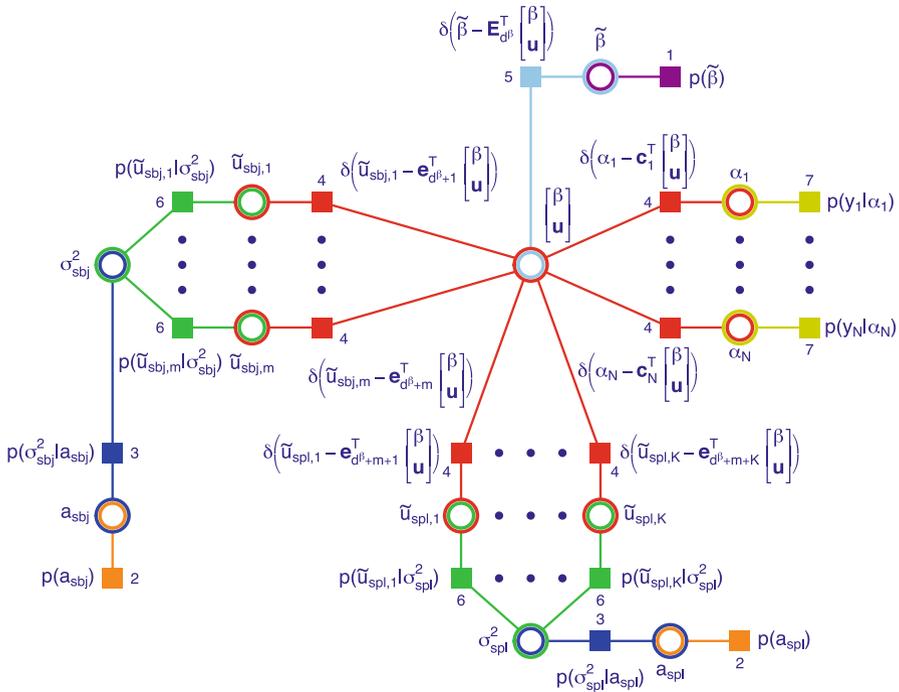


Fig. 6 Derived variable factor graph corresponding to the representation of the joint density function of random variables in the generalized additive mixed model (22) given by (23)

The posterior density function comparisons are accompanied by accuracy percentages. For a generic parameter θ , the accuracy of the approximation $q(\theta)$ to the posterior density function $p(\theta|\mathbf{y})$ is given by

$$\text{accuracy} \equiv 100 \left\{ 1 - \frac{1}{2} \int_{-\infty}^{\infty} |q(\theta) - p(\theta|\mathbf{y})| d\theta \right\} \%. \tag{24}$$

The Markov chain Monte Carlo-based posterior density functions, as well as the accuracy percentages on which they depend, are binned kernel density estimates obtained using the R function `bkde()` in the package `KernSmooth` (Wand and Ripley 2015) with direct plug-in bandwidth selection via the function `dpik()`. The density estimates should be very close to the actual posterior density functions since they are based on one million posterior draws.

We see from Fig. 7 that expectation propagation achieves excellent accuracy for the fixed effect parameters, in keeping with the simulation studies of Kim and Wand (2018). The variance parameter posterior density estimates are not as good for this particular example with accuracy scores of about 72% and 83%. Such mediocre accuracy was not apparent in the Kim and Wand (2018) simulations although their Figures 9 and 11 show accuracies for variance parameters being substantially lower than that those for fixed effect parameters. We ran the code that produced Fig. 7 on some simu-

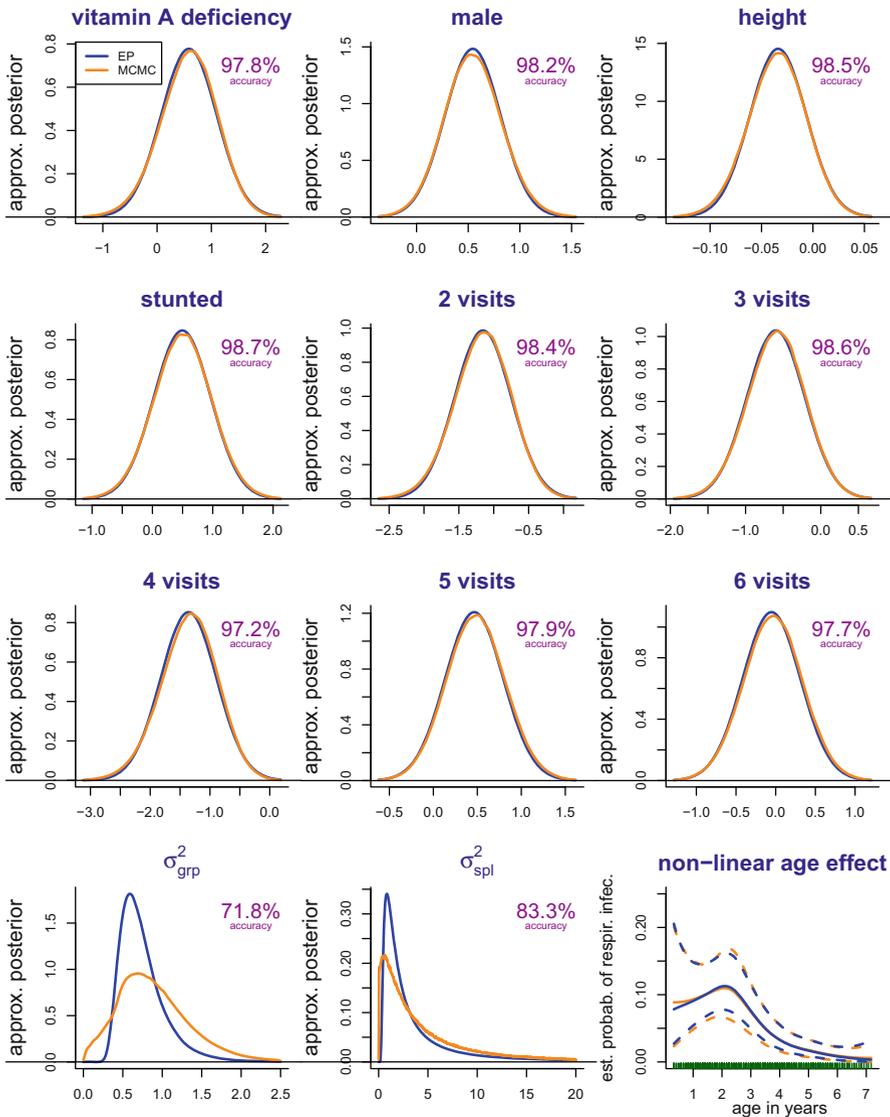
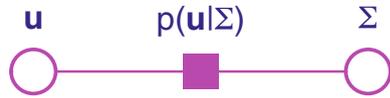


Fig. 7 Comparison of two approximate Bayesian inference methods, expectation propagation and Markov chain Monte Carlo, for model (22) applied to the Indonesian Children’s Health Study Data. The first three rows compares approximate posterior density functions for the fixed effects parameters. The heading at the top of the panel is the corresponding predictor. The first two panels in the fourth row compares approximate posterior density functions for the two variance parameters. The accuracy percentages correspond to the definition at (24). The bottom right panel compares the low-rank smoothing spline fits for the non-linear age effect on the probability of respiratory infection with all other predictors set to their averages. In this panel, the dashed curves indicate pointwise 95% credible intervals and the tick marks show the age data

Fig. 8 The factor graph fragment corresponding the factor $p(\mathbf{u}|\Sigma)$



lated data and got accuracy scores in the 85–95% range for the variance parameters. A likely explanation for the lower accuracy for variance parameters is that the posterior correlation within the pairs $(\sigma_{\text{grp}}^2, a_{\text{grp}})$ and $(\sigma_{\text{spl}}^2, a_{\text{spl}})$ is being ignored in the mean field restriction being imposed by expectation propagation in this example.

5 More elaborate expectation propagation fragments

The fragments listed in Table 2 and covered in Sect. 3 are the most fundamental ones for generalized, linear and mixed models. Whilst these fragments support expectation propagation fitting of a wide range of models, additional fragments are needed for various elaborations. We now illustrate this fact by investigating fragments needed for (a) the extension to multivariate random effects, and (b) models where the response variable is modeled according to the t distribution. As we will see, expectation propagation is quite numerically challenging for such extensions.

5.1 Multivariate random effects

The fragments in Table 2 can handle the univariate random effects structure

$$u|\sigma^2 \sim N(0, \sigma^2)$$

but they do not cover the *multivariate* random effects extension:

$$\mathbf{u}|\Sigma \sim N(\mathbf{0}, \Sigma)$$

where Σ is a unstructured $d^u \times d^u$ matrix.

The fragment corresponding to the factor

$$p(\mathbf{u}|\Sigma) = (2\pi)^{-d^u/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2} \mathbf{u}^T \Sigma^{-1} \mathbf{u})$$

is shown in Fig. 8.

Under the usual conjugacy constraints, the message from $p(\mathbf{u}|\Sigma)$ to Σ is

$$m_{p(\mathbf{u}|\Sigma) \rightarrow \Sigma}(\Sigma) = \frac{\text{proj}_{\text{IW}} [m_{\Sigma \rightarrow p(\mathbf{u}|\Sigma)}(\Sigma) \int_{\mathbb{R}^{d^u}} p(\mathbf{u}|\Sigma) m_{\mathbf{u} \rightarrow p(\mathbf{u}|\Sigma)}(\mathbf{u}) d\mathbf{u} / Z]}{m_{\Sigma \rightarrow p(\mathbf{u}|\Sigma)}(\Sigma)} \tag{25}$$

where proj_{IW} denotes projection onto the d^u -dimensional Inverse Wishart family. The messages on the right-hand side of (25) have the form

$$m_{\Sigma \rightarrow p(\mathbf{u}|\Sigma)}(\Sigma) = \exp \left\{ \left[\begin{array}{c} \log |\Sigma| \\ \text{vec}(\Sigma^{-1}) \end{array} \right]^T \boldsymbol{\eta}_{\Sigma \rightarrow p(\mathbf{u}|\Sigma)} \right\}$$

and

$$m_{\mathbf{u} \rightarrow p(\mathbf{u}|\Sigma)}(\mathbf{u}) = \exp \left\{ \left[\begin{array}{c} \mathbf{u} \\ \text{vec}(\mathbf{u}\mathbf{u}^T) \end{array} \right]^T \boldsymbol{\eta}_{\mathbf{u} \rightarrow p(\mathbf{u}|\Sigma)} \right\}$$

Introducing the shorthand

$$\boldsymbol{\eta}^\heartsuit \equiv \boldsymbol{\eta}_{\Sigma \rightarrow p(\mathbf{u}|\Sigma)} \quad \text{and} \quad \boldsymbol{\eta}^\diamond \equiv \boldsymbol{\eta}_{\mathbf{u} \rightarrow p(\mathbf{u}|\Sigma)},$$

arguments analogous to those given in Section S.2.5 of the online supplement lead to the function of Σ inside the $\text{proj}_{\text{IW}}[\cdot]$ in (25) being proportional to

$$|\Sigma|^{\eta_1^\heartsuit} |2\Sigma \text{vec}^{-1}(\boldsymbol{\eta}_2^\diamond) - \mathbf{I}|^{-1/2} \text{tr}\{\Sigma^{-1} \text{vec}^{-1}(\boldsymbol{\eta}_2^\heartsuit)\} \times \exp \left[-\frac{1}{4}(\boldsymbol{\eta}_1^\diamond)^T \{2\Sigma \text{vec}^{-1}(\boldsymbol{\eta}_2^\diamond) - \mathbf{I}\} \{\text{vec}^{-1}(\boldsymbol{\eta}_2^\diamond)\}^{-1} \boldsymbol{\eta}_1^\diamond \right]. \tag{26}$$

The next step is to compute $E\{\log |\Sigma|\}$ and $E(\Sigma^{-1})$ with expectation with respect to the density function obtained by normalizing (26). This is a particularly challenging numerical problem since it involves numerical integration of the cone of $d^u \times d^u$ symmetric positive definite matrices. Then

$$\boldsymbol{\eta}_{p(\mathbf{u}|\Sigma) \rightarrow \Sigma} = (\nabla A)_{\text{IW}}^{-1} \left(\left[\begin{array}{c} E\{\log |\Sigma|\} \\ E\{\text{vech}(\Sigma^{-1})\} \end{array} \right] \right). \tag{27}$$

Note that the function $(\nabla A)_{\text{IW}}$ admits the explicit form

$$(\nabla A)_{\text{IW}} \left(\left[\begin{array}{c} \eta_1 \\ \boldsymbol{\eta}_2 \end{array} \right] \right) = \left[\begin{array}{c} \log \left| -\text{vech}^{-1}(\boldsymbol{\eta}_2) \right| - \sum_{j=1}^{d^u} \text{digamma} \left(-\eta_1 - \frac{1}{2}(d^u + j) \right) \\ \left\{ \eta_1 + \frac{1}{2}(d^u + 1) \right\} \text{vech}[\{\text{vech}^{-1}(\boldsymbol{\eta}_2)\}^{-1}] \end{array} \right]$$

where $[\eta_1 \ \boldsymbol{\eta}_2^T]^T$ is the partition of the natural parameter vector into the first entry (η_1) and the remaining $\frac{1}{2}d^u(d^u + 1)$ entries ($\boldsymbol{\eta}_2$). However, evaluation of (27) involves numerical inversion of $(\nabla A)_{\text{IW}}$ in $\{1 + \frac{1}{2}d^u(d^u + 1)\}$ -dimensional space.

In conclusion, literal application of expectation propagation for multivariate random effects is quite daunting and effective implementation for even $2 \leq d^u \leq 5$ is a very challenging numerical problem.

5.2 *t* likelihood

The Gaussian fragment, treated in Sect. 3.6, corresponds to the specification $y|\alpha, \sigma^2 \sim N(\alpha, \sigma^2)$. Now consider the extension to the *t* distribution:

$$y|\alpha, \sigma^2, \nu \sim t(\alpha, \sigma^2, \nu) \tag{28}$$

where $\nu > 0$ is the degrees of freedom parameter. Low values of ν correspond to heavy-tailed distributions. The Gaussian likelihood is the $\nu \rightarrow \infty$ limiting case. The density function corresponding to (28) is

$$p(y|\alpha, \sigma^2, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu}\Gamma(\nu/2)\{1 + (y - \alpha)^2/(\nu\sigma^2)\}^{\frac{\nu+1}{2}}}.$$

One could work with this density function in the expectation propagation message Eqs. (8) and (9), but trivariate numerical integration is required. In other Bayesian computation contexts such as Markov chain Monte Carlo (e.g. Verdinelli and Wasserman 1991) and variational message passing (e.g. McLean and Wand 2018) it is common to replace (28) by the auxiliary variable representation

$$y|\alpha, \sigma^2, a \sim N(\alpha, a\sigma^2), \quad a|\nu \sim \text{Inverse-}\chi^2(\nu, \nu) \tag{29}$$

to aid tractability. Expectation propagation also benefits from this representation of the *t*-likelihood specification. The fragments corresponding to the factor product

$$p(y|\alpha, \sigma^2, a) p(a|\nu) p(\nu) \tag{30}$$

are shown in Fig. 9.

None of these fragments are among those treated in Sect. 3. Therefore, extension to *t* likelihood models requires expectation propagation updates for these three new fragments. Unfortunately, as we will see, difficult numerical challenges arise for these updates. We now focus on each one in turn.

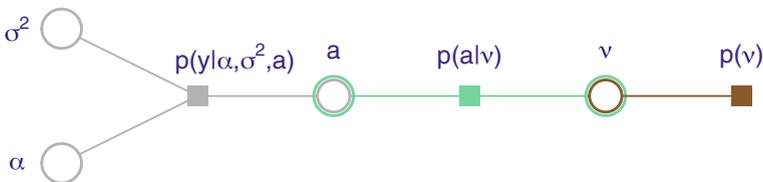


Fig. 9 Color-coded fragments corresponding to the factors $p(y|\alpha, \sigma^2, a)$, $p(a|\nu)$ and $p(\nu)$ appearing in (30)

5.2.1 The $p(y|\alpha, \sigma^2, a)$ fragment

The factor for this fragment is

$$p(y|\alpha, \sigma^2, a) = (2\pi a\sigma^2)^{-1/2} \exp\left\{-\frac{(y - \alpha)^2}{2a\sigma^2}\right\}.$$

Conjugacy considerations dictate the assumption:

(31)

all messages passed to α from factors outside of the fragment are in the Univariate Normal family and all messages passed to either a or σ^2 from factors outside of the fragment are in the Inverse Chi-Squared family.

This leads to the factor to stochastic node messages taking the forms:

$$m_{p(y|\alpha, \sigma^2, a) \rightarrow \alpha}(\alpha) = \exp\left\{\left[\begin{matrix} \alpha \\ \alpha^2 \end{matrix}\right]^T \eta_{p(y|\alpha, \sigma^2, a) \rightarrow \alpha}\right\},$$

$$m_{p(y|\alpha, \sigma^2, a) \rightarrow \sigma^2}(\sigma^2) = \exp\left\{\left[\begin{matrix} \log(\sigma^2) \\ 1/\sigma^2 \end{matrix}\right]^T \eta_{p(y|\alpha, \sigma^2, a) \rightarrow \sigma^2}\right\}$$

and

$$m_{p(y|\alpha, \sigma^2, a) \rightarrow a}(a) = \exp\left\{\left[\begin{matrix} \log(a) \\ 1/a \end{matrix}\right]^T \eta_{p(y|\alpha, \sigma^2, a) \rightarrow a}\right\}.$$

The derivations of the natural parameter updates are similar in nature to those given in Section S.2.6 of the online supplement for the Gaussian fragment. However, the form $a\sigma^2$ (rather than σ^2) in the variance means that the natural parameter updates require evaluation of the bivariate integral-defined function

$$\mathcal{B}_2(p, q_1, q_2, r_1, r_2, s, t, u) \equiv \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{x_1^p \exp\{q_1x_1 + q_2x_2 - r_1e^{x_1} - r_2e^{x_2} - se^{x_1+x_2}/(t + e^{x_1+x_2})\}}{(t + e^{x_1+x_2})^u} dx_1 dx_2$$

for

$$p \geq 0, q_1, q_2 \in \mathbb{R}, r_1, r_2 > 0, s \geq 0, t > 0, u > 0$$

rather than the univariate integral-defined function $\mathcal{B}(p, q, r, s, t, u)$ given by equations (S.1) of the online supplement.

5.2.2 The $p(a|v)$ fragment

The relevant factor is

$$p(a|v) = \frac{(v/2)^{v/2}}{\Gamma(v/2)} a^{-(v/2)-1} \exp\{-(v/2)/a\}, \quad a, v > 0.$$

Let $v \equiv v/2$ be a simple linear transformation of v . For the remainder of this section we work with v , rather than v , since it leads to a simpler exposition. Now note that

$$p(a|v) \propto \begin{cases} \exp \left\{ \begin{bmatrix} \log(a) \\ 1/a \end{bmatrix}^T \begin{bmatrix} -v \\ -v-1 \end{bmatrix} \right\} & \text{as a function of } a, \\ \exp \left\{ \begin{bmatrix} v \log(v) - \log\{\Gamma(v)\} \\ v \end{bmatrix}^T \begin{bmatrix} 1 \\ -1/a - \log(a) \end{bmatrix} \right\} & \text{as a function of } v. \end{cases}$$

To ensure conjugacy we should then impose the restriction:

$$\begin{aligned} &\text{all messages passed to } a \text{ from factors outside of the} \\ &\text{fragment are in the Inverse Chi-Squared family and all} \\ &\text{messages passed to either } v \text{ from factors outside} \\ &\text{of the fragment are in the Moon Rock family.} \end{aligned} \tag{32}$$

The definition of the Moon Rock family is given in Table 1. The messages passed from $p(a|v)$ are then of the form

$$m_{p(a|v) \rightarrow a}(a) = \exp \left\{ \begin{bmatrix} \log(a) \\ 1/a \end{bmatrix}^T \eta_{p(a|v) \rightarrow a} \right\}$$

and

$$m_{p(a|v) \rightarrow v}(v) = \exp \left\{ \begin{bmatrix} v \log(v) - \log\{\Gamma(v)\} \\ v \end{bmatrix}^T \eta_{p(a|v) \rightarrow v} \right\}.$$

The message $m_{p(a|v) \rightarrow a}(a)$ has a treatment similar to that for $m_{p(\sigma^2|a) \rightarrow \sigma^2}(\sigma^2)$ and $m_{p(\sigma^2|a) \rightarrow a}(a)$ in Section S.2.3 of the online supplement and $m_{p(by|\alpha, \sigma^2) \rightarrow \sigma^2}(\sigma^2)$ in Section S.2.6 of the online supplement, with projection onto the Inverse Chi-Squared family, although bivariate numerical integration is required. On the other hand,

$$m_{p(a|v) \rightarrow v}(v) = \frac{\text{proj}_{\text{MR}} \left[m_{v \rightarrow p(a|v)}(v) \int_0^\infty p(a|v) m_{p(a|v) \rightarrow a}(a) da / Z \right]}{m_{v \rightarrow p(a|v)}(v)}.$$

where proj_{MR} denotes projection onto the Moon Rock family. The function of ν inside the $\text{proj}_{\text{MR}}[\]$ is proportional to

$$h(\nu) \equiv \{\nu^\nu / \Gamma(\nu)\}^{\eta_1^\sharp+1} e^{\eta_2^\flat \nu} \Gamma(\nu - \eta_1^\sharp) / (\nu + \eta_2^\flat)^{\nu - \eta_1^\sharp}$$

where

$$\eta^\sharp \equiv \eta_{p(a|\nu) \rightarrow a} \quad \text{and} \quad \eta^\flat \equiv \eta_{\nu \rightarrow p(a|\nu)}.$$

Then

$$\eta_{p(a|\nu) \rightarrow \nu} = (\nabla A_{\text{MR}})^{-1} \left(\begin{bmatrix} \int_0^\infty \{\nu \log(\nu) - \log \Gamma(\nu)\} h(\nu) d\nu / \int_0^\infty h(\nu) d\nu \\ \int_0^\infty \nu h(\nu) d\nu / \int_0^\infty h(\nu) d\nu \end{bmatrix} \right)$$

where

$$A_{\text{MR}}(\eta) \equiv \log \left[\int_0^\infty \{t^t / \Gamma(t)\}^{\eta_1} \exp(\eta_2 t) dt \right]$$

is the log-partition function of the Moon Rock exponential family. This implies that

$$(\nabla A_{\text{MR}})(\eta) = \begin{bmatrix} \int_0^\infty \{t \log(t) - \log \Gamma(t)\} \{t^t / \Gamma(t)\}^{\eta_1} \exp(\eta_2 t) dt / \int_0^\infty \{t^t / \Gamma(t)\}^{\eta_1} \exp(\eta_2 t) dt \\ \int_0^\infty t \{t^t / \Gamma(t)\}^{\eta_1} \exp(\eta_2 t) dt / \int_0^\infty \{t^t / \Gamma(t)\}^{\eta_1} \exp(\eta_2 t) dt \end{bmatrix}.$$

This particular exponential family is not well-studied and we are not aware of any published theory concerning the properties of ∇A_{MR} and $(\nabla A_{\text{MR}})^{-1}$. Standard analytic arguments can be used to show that the domain of ∇A_{MR} is

$$H = \left\{ \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} : \eta_1 \geq 0, \eta_1 + \eta_2 < 0 \right\}.$$

It is conjectured that the image of H under ∇A_{MR} is

$$T = \left\{ \begin{bmatrix} \tau_1 \\ \tau_2 \end{bmatrix} : \tau_1 < \tau_2 \log(\tau_2) - \log \Gamma(\tau_2) \right\}.$$

Figure 10 shows the domain of ∇A_{MR} and the conjectured domain of $(\nabla A_{\text{MR}})^{-1}$ as well as some example mappings between the two spaces.

Evaluation of $(\nabla A_{\text{MR}})^{-1}$ is a non-trivial problem. It requires numerical inversion techniques such as Newton-Raphson iteration. Moreover, each of the iterative updates involves evaluation of ∇A_{MR} and, possibly, its first partial derivatives. None of these functions are available in closed form and require numerical integration.

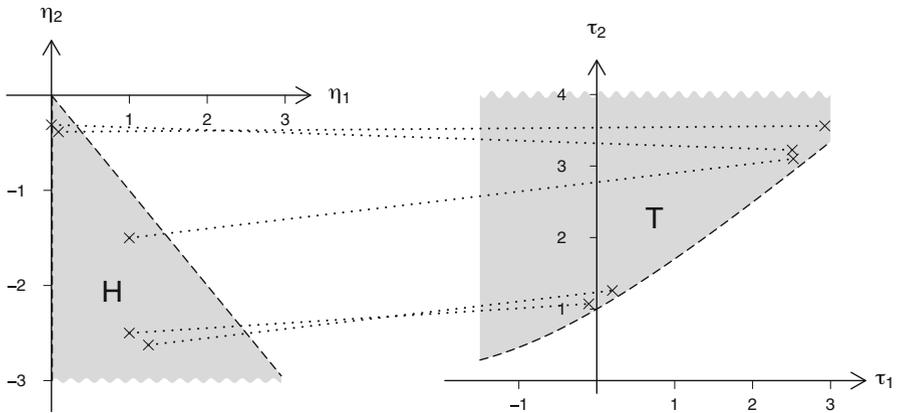


Fig. 10 Illustration of the bijective maps between H and T for the Moon Rock exponential family. The crosses and dotted lines depict five example $\eta \in H$ and $\tau = \nabla A_{MR}(\eta) \in T$ pairs. Since ∇A_{MR} is a bijective map, the crosses and dotted lines equivalently depict five example $\tau \in T$ and $\eta = (\nabla A_{MR})^{-1}(\tau) \in H$ pairs

5.2.3 The $p(v)$ fragment

This simple fragment has the factor to stochastic node message

$$m_{p(v) \rightarrow v}(v) \propto p(v)$$

corresponding to the prior distribution on v . The conjugate family of prior density functions is

$$p(v) \propto \{(v/2)^{v/2} / \Gamma(v/2)\}^{A_v} \exp(-\frac{1}{2} B_v v), \quad v > 0.$$

for hyperparameters $A_v \geq 0$ and $B_v > A_v$.

In terms of $v = v/2$, the relevant message is

$$m_{p(v) \rightarrow v}(v) = \exp \left\{ \begin{bmatrix} v \log(v) - \log\{\Gamma(v)\} \\ v \end{bmatrix}^T \begin{bmatrix} A_v \\ -B_v \end{bmatrix} \right\}.$$

5.3 Summary of numerical challenges

The previous two subsections make it clear that elaborations such as multivariate random effects and fancier likelihoods involve profound numerical challenges for the expectation propagation paradigm. Table 3 summarizes the numerical challenges of all of the non-trivial fragments treated in this article.

The first ten fragments in Table 3 have the attraction of requiring only numerical evaluation of univariate integral within the families given by equations (S.1) and (S.2) of the online supplement. The probit likelihood fragment stands out as a special case of

Table 3 The numerical integration and Kullback–Leibler projection demands of the non-trivial fragments discussed in this article

Fragment name	Numeric. integrat. demands	Kull.-Leib. projec. demands
Gaussian prior	None	None
Inverse Wishart prior	None	None
Moon Rock prior	None	None
Iterated Inverse Chi Squared	Univariate quadrature	Inversion of log –digamma
Linear comb. deriv. var.	None	None
Multiv. lin. comb. deriv. var.	None	None
Gaussian	Univariate quadrature	Inversion of log –digamma
Logistic likelihood	Univariate quadrature	Trivial
Probit likelihood	None	Trivial
Poisson likelihood	Univariate quadrature	Trivial
Multiple random effects	Multivariate quadrature	Inversion of a multivariate function
t likelihood direct	Trivariate quadrature	Inversion of log –digamma and a non-explicit bivariate function
t likelihood aux. var.	Bivariate quadrature	Inversion of log –digamma and a non-explicit bivariate function

a likelihood that does not require any numerical methods for expectation propagation message passing.

The last three fragments of Table 3 are considerably more demanding in terms of numerical analysis and exact expectation propagation may not be viable. One option is to call upon Monte Carlo methodology and approximate the messages stochastically as has been proposed in the more recent expectation propagation literature such as Heese et al. (2013) and Lienart et al. (2015). Whilst such Monte Carlo approximation represents a departure from the fast deterministic approximate inference aspect of expectation propagation and mean field variational Bayes, it may be the only practical way by which such fragments can be handled.

Acknowledgements This research was supported by Australian Research Council Discovery Project DP180100597.

References

- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., et al. (2017). Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1), 1–32.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2014). *Bayesian data analysis* (3rd ed.). Boca Raton: CRC Press.
- Guo, J., Gabry, J. & Goodrich, B. (2017). The R package rstan: R interface to Stan. R package (version 2.17.2). <http://mc-stan.org>.
- Heese, N., Tarlow, D., & Winn, J. (2013). Learning to pass expectation propagation messages. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems* (Vol. 26, pp. 3219–3227). Red Hook: Curran Associates, Incorporated.

- Herbrich, R. (2005). Gaussian expectation propagation. <https://www.microsoft.com/en-us/research/publication/on-gaussian-expectation-propagation/>.
- Heskes, T., Opper, M., Wiegerinck, W., Winther, O., & Zoeter, O. (2005). Approximate inference techniques with expectation constraints. *Journal of Statistical Mechanics Theory and Experiment*, *P11015*, 1–24.
- Heskes, T., & Zoeter, O. (2002). Expectation propagation for approximate inference in dynamic Bayesian networks. In A. Darwiche & N. Friedman (Eds.), *Proceedings of the eighteenth annual conference on uncertainty in artificial intelligence* (pp. 216–223). San Francisco: Morgan Kaufmann.
- Jylänki, P., Vanhatalo, J., & Vehtari, A. (2011). Robust Gaussian process regression with a student-*t* likelihood. *Journal of Machine Learning Research*, *12*, 3227–3257.
- Kim, A. S. I., & Wand, M. P. (2016). The explicit form of expectation propagation for a simple statistical model. *Electronic Journal of Statistics*, *10*, 550–581.
- Kim, A. S. I., & Wand, M. P. (2018). On expectation propagation for generalized, linear and mixed models. *Australian and New Zealand Journal of Statistics*, *60*, 75–102.
- Lienart, T., Teh, Y. W., & Doucet, A. (2015). Expectation particle belief propagation. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 28, pp. 3609–3617). Red Hook: Curran Associates, Incorporated.
- McLean, M. W., & Wand, M. P. (2018). Variational message passing for elaborate response regression models. *Bayesian Analysis* (in press).
- Minka, T. (2005). Divergence measures and message passing. *Microsoft research technical report series*, MSR-TR-2005-173, pp. 1–17.
- Minka, T., & Winn, J. (2008). Gates: A graphical notation for mixture models. *Microsoft research technical report series*, MSR-TR-2008-185, pp. 1–16.
- Minka, T. P. (2001). Expectation propagation for approximate Bayesian inference. In J. S. Breese & D. Koller (Eds.), *Proceedings of the seventeenth conference on uncertainty in artificial intelligence* (pp. 362–369). Burlington: Morgan Kaufmann.
- Minka, T., Winn, J. M., Guiver, J. P., Webster, S., Zaykov, Y., Yangel, B., Spengler, A. & Bronskill, J. (2014). Infer.NET 2.6, Microsoft Research Cambridge, 2014. <http://research.microsoft.com/infernet>.
- Murphy, K. (2007). Software for graphical models: a review. *International Society for Bayesian Analysis Bulletin*, *14*, 13–15.
- Nolan, T. H., & Wand, M. P. (2017). Accurate logistic variational message passing: Algebraic and numerical details. *Stat*, *6*, 102–112.
- Opper, M., & Winther, O. (2000). Gaussian processes for classification: Mean-field algorithms. *Neural Computation*, *12*, 2655–2684.
- Opper, M., & Winther, O. (2005). Expectation consistent approximate inference. *Journal of Machine Learning Research*, *6*, 2177–2204.
- Ruppert, D., Wand, M. P., & Carroll, R. J. (2009). Semiparametric regression during 2003–2007. *Electronic Journal of Statistics*, *3*, 1193–1256.
- Sommer, A. (1982). *Nutritional blindness*. New York: Oxford University Press.
- Thouless, D. J., Anderson, P. W., & Palmer, R. G. (1977). Solution of a “solvable model of a spin glass”. *The Philosophical Magazine*, *35*, 593.
- Verdinelli, I., & Wasserman, L. (1991). Bayesian analysis of outlier problems using the Gibbs sampler. *Statistics and Computing*, *1*(2), 105–117.
- Wainwright, M. J., & Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, *1*, 1–305.
- 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.
- Wand, M. P., & Ripley, B. D. (2015). The R package `kernelSmooth`. Functions for kernel smoothing supporting Wand & Jones (1995) (version 2.23). <https://cran.R-project.org>.
- Wand, M. P., Ormerod, J. T., Padoan, S. A., & Frühwirth, R. (2011). Mean Field Variational Bayes for Elaborate Distributions. *Bayesian Analysis*, *6*(4), 847–900.