



The Inverse G-Wishart distribution and variational message passing

Luca Maestrini* and Matt P. Wand

University of Technology Sydney

Summary

Message passing on a factor graph is a powerful paradigm for the coding of approximate inference algorithms for arbitrarily large graphical models. The notion of a factor graph fragment allows for compartmentalisation of algebra and computer code. We show that the Inverse G-Wishart family of distributions enables fundamental variational message passing factor graph fragments to be expressed elegantly and succinctly. Such fragments arise in models for which approximate inference concerning covariance matrix or variance parameters is made, and are ubiquitous in contemporary statistics and machine learning.

Key words: approximate Bayesian inference; G-Wishart distribution; mean-field variational Bayes; scalable statistical methodology

1. Introduction

We argue that a very general family of covariance matrix distributions, known as the *Inverse G-Wishart* family, plays a fundamental role in modularisation of variational inference algorithms via variational message passing when a factor graph fragment (Wand 2017) approach is used. A factor graph fragment, or *fragment* for short, is a sub-graph of the relevant factor graph consisting of a factor and all of its neighbouring nodes. Even though use of the Inverse G-Wishart distribution is not necessary, its adoption allows for fundamental factor graph fragment natural parameter updates to be expressed elegantly and succinctly. An essential aspect of this strategy is that the Inverse G-Wishart distribution is the *only* distribution used for covariance matrix and variance parameters. The family includes as special cases the Inverse Chi-Squared, Inverse Gamma and Inverse Wishart distributions. Therefore, just a single distribution is required which leads to savings in notation and code. While similar comments concerning modularity apply to Monte Carlo-based approaches to approximate Bayesian inference, here we focus on variational inference.

Two of the most common contemporary approaches to fast approximate Bayesian inference are mean-field variational Bayes (e.g. Attias 1999) and expectation propagation (e.g. Minka 2001). Minka (2005) explains how each approach can be expressed as message passing on relevant *factor graphs* with *variational message passing* (Winn & Bishop 2005) being the name used for the message passing version of mean-field variational Bayes. Wand (2017) introduces the concept of *factor graph fragments*, or *fragments* for short, for compartmentalisation of variational message passing into atom-like components. Chen &

*Author to whom correspondence should be addressed.

School of Mathematical and Physical Sciences, University of Technology Sydney, P.O. Box 123, Broadway, NSW 2007, Australia. email: luca.maestrini@uts.edu.au.

Wand (2020) demonstrate the use of fragments for expectation propagation. Explanations of factor graph-based variational message passing that match the current exposition are given in Wand (2017 Sections 2.4–2.5).

Wand (2017 Sections 4.1.2–4.1.3) introduces two variational message passing fragments known as the *Inverse Wishart prior* fragment and the *iterated Inverse G-Wishart* fragment. The first of these simply corresponds to imposing an Inverse Wishart prior on a covariance matrix. In the scalar case, this reduces to imposing an Inverse Chi-Squared or, equivalently, an Inverse Gamma prior on a variance parameter. The iterated Inverse G-Wishart fragment facilitates the imposition of arbitrarily non-informative priors on standard deviation parameters such as members of the Half- t family (Gelman 2006; Polson & Scott 2012). An extension to the covariance matrix case, for which there is the option to impose marginal Uniform distribution priors over the interval $(-1, 1)$ on correlation parameters, is elucidated in Huang & Wand (2013). Mulder & Pericchi (2018) provide a different type of extension that is labelled the *Matrix-F* distribution. These two fragments arise in many classes of Bayesian models, such as both Gaussian and generalised response linear mixed models (e.g. McCulloch, Searle & Neuhaus 2008), Bayesian factor models (e.g. Conti *et al.* 2014), vector autoregressive models (e.g. Assaf *et al.* 2019), and generalised additive mixed models and group-specific curve models (e.g. Harezlak, Ruppert & Wand 2018).

Despite the fundamentalness of Inverse G-Wishart-based fragments for variational message passing, the main reference to date, Wand (2017), is brief in its exposition and contains some errors that affect certain cases. In this article, we provide a detailed exposition of the Inverse G-Wishart distribution in the context of variational message passing and list the Inverse Wishart prior and iterated Inverse G-Wishart fragment updates in full ready-to-code forms. R functions (R Core Team 2021) that implement these algorithms are provided as part of the supplementary material of this article. We also explain the errors in Wand (2017).

Section 2 contains relevant definitions and results concerning the G-Wishart and Inverse G-Wishart distributions. Connections with the Huang–Wand and *Matrix-F* families of marginally noninformative prior distributions for covariance matrices are summarized in Section 3 and in Section 4 we point to background material on variational message passing. In Sections 5 and 6, we provide detailed accounts of the two variational message passing fragments pertaining to variance and covariance matrix parameters, expanding on what is presented in Wand (2017 Sections 4.1.2 and 4.1.3), and making some corrections to what is presented there. In Section 7, we provide explicit instructions on how the two fragments are used to specify different types of prior distributions on standard deviation and covariance matrix parameters in variational message passing-based approximate Bayesian inference. Section 8 contains an data analytic example that illustrates the use of the covariance matrix fragment update algorithms. Some closing discussion is given in Section 9. A web-supplement contains relevant details.

2. The G-Wishart and Inverse G-Wishart distributions

A random matrix X has an Inverse G-Wishart distribution if and only if X^{-1} has a G-Wishart distribution. In this section, we first review the G-Wishart distribution, which has an established literature. Then we discuss the Inverse G-Wishart distribution and list properties that are relevant to its employment in variational message passing.

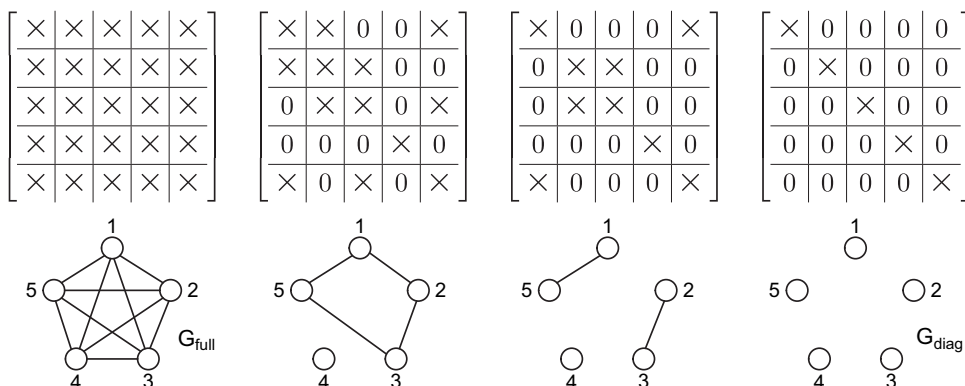


Figure 1. The zero/non-zero entries of four \$5 \times 5\$ symmetric matrices with non-zero entries denoted by \$\times\$. Underneath each matrix is the 5-node undirected graph that the matrix respects. The nodes are numbered according to the rows and columns of the matrices. A graph edge is present between nodes \$i\$ and \$j\$ whenever the \$(i,j)\$ entry of the matrix is non-zero. The graph respected by the full matrix is denoted by \$G_{full}\$. The graph respected by the diagonal matrix is denoted by \$G_{diag}\$.

Let \$G\$ be an undirected graph with \$d\$ nodes labelled \$1, \dots, d\$ and set \$E\$ consisting of pairs of nodes that are connected by an edge. We say that the symmetric \$d \times d\$ matrix \$M\$ respects \$G\$ if

$$M_{ij} = 0, \quad \text{for all } \{i, j\} \notin E.$$

Figure 1 shows the zero/non-zero entries of four \$5 \times 5\$ symmetric matrices. For each matrix, the 5-node graph that the matrix respects is shown underneath.

The first graph in Figure 1 is totally connected and corresponds to the matrix being full. Hence, we denote this graph by \$G_{full}\$. At the other end of the spectrum is the last graph of Figure 1, which is totally disconnected. Since this corresponds to the matrix being diagonal, we denote this graph by \$G_{diag}\$.

An important concept in G-Wishart and Inverse G-Wishart distribution theory is graph decomposability. An undirected graph \$G\$ is *decomposable* if and only if all cycles of four or more nodes have an edge that is not part of the cycle but connects two nodes of the cycle. In Figure 1, the first, third and fourth graphs are decomposable. However, the second graph is not decomposable since it contains a four-node cycle that is devoid of edges that connect pairs of nodes within this cycle. Alternative labels for decomposable graphs are *chordal* graphs and *triangulated* graphs.

In Sections 2.1 and 2.2 we define the G-Wishart and Inverse G-Wishart distributions and treat important special cases. This exposition depends on particular notation, which we define here. For a generic proposition \$\mathcal{P}\$, we define \$\mathbf{1}(\mathcal{P})\$ to equal 1 if \$\mathcal{P}\$ is true and zero otherwise. If the random variables \$x_j, 1 \leq j \leq d\$, are independent such that \$x_j\$ has distribution \$\mathcal{D}_j\$ we write \$x_j \overset{\text{ind.}}{\sim} \mathcal{D}_j, 1 \leq j \leq d\$. For a \$d \times 1\$ vector \$\mathbf{v}\$, let \$\text{diag}(\mathbf{v})\$ be the \$d \times d\$ diagonal matrix with diagonal comprising the entries of \$\mathbf{v}\$ in order. For a \$d \times d\$ matrix \$\mathbf{M}\$, let \$\text{diagonal}(\mathbf{M})\$ denote the \$d \times 1\$ vector comprising the diagonal entries of \$\mathbf{M}\$ in order. The \$\text{vec}\$ and \$\text{vech}\$ matrix operators are well-established (e.g. Gentle 2007). If \$\mathbf{a}\$ is a \$d^2 \times 1\$ vector

then $\text{vec}^{-1}(\mathbf{a})$ is the $d \times d$ matrix such that $\text{vec}(\text{vec}^{-1}(\mathbf{a})) = \mathbf{a}$. The matrix \mathbf{D}_d , known as the *duplication matrix of order d* , is the $d^2 \times \{d(d+1)/2\}$ matrix containing only zeros and ones such that $\mathbf{D}_d \text{vech}(\mathbf{A}) = \text{vec}(\mathbf{A})$ for any symmetric $d \times d$ matrix \mathbf{A} (Magnus & Neudecker 2019). For example,

$$\mathbf{D}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

The Moore–Penrose inverse of \mathbf{D}_d is $\mathbf{D}_d^+ = (\mathbf{D}_d^\top \mathbf{D}_d)^{-1} \mathbf{D}_d^\top$ and is such that $\mathbf{D}_d^+ \text{vec}(\mathbf{A}) = \text{vech}(\mathbf{A})$ for a symmetric matrix \mathbf{A} .

2.1. The G-Wishart distribution

The *G-Wishart distribution* (Atay-Kayis & Massam 2005) is defined as follows:

Definition 1 Let \mathbf{X} be a $d \times d$ symmetric and positive definite random matrix and G be a d -node undirected graph such that \mathbf{X} respects G . For $\delta > 0$ and a symmetric positive definite $d \times d$ matrix $\mathbf{\Lambda}$ we say that \mathbf{X} has a G-Wishart distribution with graph G , shape parameter δ and rate matrix $\mathbf{\Lambda}$, and write

$$\mathbf{X} \sim \text{G-Wishart}(G, \delta, \mathbf{\Lambda}),$$

if and only if the non-zero values of the density function of \mathbf{X} satisfy

$$p(\mathbf{X}) \propto |\mathbf{X}|^{(\delta-2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda X}) \right\}. \tag{1}$$

Obtaining an expression for the normalising factor of a general G-Wishart density function is a challenging problem and recently was resolved by Uhler, Lenkoski & Richards (2018). In the special case where G is a decomposable graph, a relatively simple expression for the normalising factor exists and is given, for example, by (1.4) of Uhler *et al.* (2018). The non-decomposable case is much more difficult and treated in Section 3 of Uhler *et al.* (2018), but the normalising factor does not have a succinct expression for general G . Similar comments apply to expressions for the mean of a G-Wishart random matrix. As discussed in Atay-Kayis & Massam (2005 section 3), the G-Wishart distribution has connections with other distributional constructs such as the hyper Wishart law defined by Dawid & Lauritzen (1993).

Let G_{full} be the totally connected d -node undirected graph and G_{diag} be the totally disconnected d -node undirected graph. The special cases of $G = G_{\text{full}}$ and $G = G_{\text{diag}}$ are such that the normalising factor and mean do have simple closed-form expressions. Since these cases arise in fundamental variational message passing algorithms, we now turn our attention to them.

2.1.1. The $G = G_{\text{full}}$ special case

In the case where G is a fully connected graph we have:

Result 1 If the $d \times d$ random matrix \mathbf{X} is such that $\mathbf{X} \sim \text{G-Wishart}(G_{\text{full}}, \delta, \mathbf{\Lambda})$ then

$$p(\mathbf{X}) = \frac{|\mathbf{\Lambda}|^{(\delta+d-1)/2}}{2^{d(\delta+d-1)/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma\left(\frac{\delta+d-j}{2}\right)} |\mathbf{X}|^{(\delta-2)/2} \exp\left\{-\frac{1}{2} \text{tr}(\mathbf{\Lambda X})\right\} \times \mathbf{1}(\mathbf{X} \text{ a symmetric and positive definite } d \times d \text{ matrix}). \tag{2}$$

The mean of \mathbf{X} is

$$E(\mathbf{X}) = (\delta + d - 1)\mathbf{\Lambda}^{-1}.$$

Result 1 is not novel at all since the $G = G_{\text{full}}$ case corresponds to \mathbf{X} having a Wishart distribution. In other words, (2) is simply the density function of a Wishart random matrix. However, it is worth pointing out the shape parameter used here is different from that commonly used for the Wishart distribution. For example, in Gelman *et al.* (2014 table A.1) the shape parameter is denoted by ν and is related to the shape parameter of (2) according to

$$\nu = \delta + d - 1$$

and therefore are the same only in the special case of \mathbf{X} being scalar. Also, note that Definition 1 and Result 1 use the rate matrix parameterisation, whereas Gelman *et al.* (2014) Table A.1 use the scale matrix parameterisation for the Wishart distribution. The scale matrix is $\mathbf{\Lambda}^{-1}$.

2.1.2. The $G = G_{\text{diag}}$ special case

Before treating the $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$ situation, we define the notation

$$x \sim \text{Gamma}(\alpha, \beta) \tag{3}$$

to mean that the scalar random variable x has a Gamma distribution with shape parameter α and rate parameter β . The density function corresponding to (3) is

$$p(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x) \mathbf{1}(x > 0).$$

The $\text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$ distribution is tied intimately to the Gamma distribution, as Result 2 shows.

Result 2 Suppose that the $d \times d$ random matrix \mathbf{X} is such that $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$. Then the non-zero entries of \mathbf{X} satisfy

$$X_{jj} \stackrel{\text{ind.}}{\sim} \text{Gamma}\left(\frac{1}{2}\delta, \frac{1}{2}\Lambda_{jj}\right), \quad 1 \leq j \leq d,$$

where Λ_{jj} is the j th diagonal entry of $\mathbf{\Lambda}$. The density function of \mathbf{X} is

$$\begin{aligned} p(\mathbf{X}) &= \frac{|\mathbf{\Lambda}|^{\delta/2}}{2^{d\delta/2} \Gamma(\delta/2)^d} |\mathbf{X}|^{(\delta-2)/2} \exp\left\{-\frac{1}{2} \text{tr}(\mathbf{\Lambda X})\right\} \prod_{j=1}^d \mathbf{1}(X_{jj} > 0) \\ &= \frac{\prod_{j=1}^d \Lambda_{jj}^{\delta/2}}{2^{d\delta/2} \Gamma(\delta/2)^d} \prod_{j=1}^d X_{jj}^{(\delta-2)/2} \exp\left(-\frac{1}{2} \sum_{j=1}^d \Lambda_{jj} X_{jj}\right) \prod_{j=1}^d \mathbf{1}(X_{jj} > 0). \end{aligned}$$

The mean of \mathbf{X} is

$$E(\mathbf{X}) = \delta \mathbf{\Lambda}^{-1} = \delta \text{diag}(1/\Lambda_{11}, \dots, 1/\Lambda_{dd}).$$

We now make some remarks concerning Result 2.

- (i) When $G = G_{\text{diag}}$ the off-diagonal entries of $\mathbf{\Lambda}$ have no effect on the distribution of \mathbf{X} . In other words, the declaration $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$ is equivalent to the declaration $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \text{diag}\{\text{diagonal}(\mathbf{\Lambda})\})$.
- (ii) The declaration $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$ is equivalent to the diagonal entries of \mathbf{X} being independent Gamma random variables with shape parameter $(1/2)\delta$ and rate parameters equalling the diagonal entries of $(1/2)\mathbf{\Lambda}$.
- (iii) Even though statements concerning the distributions of independent random variables may seem simpler than a statement of the form $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \mathbf{\Lambda})$, the major thrust of this article is the elegance provided by key variational message passing fragment updates being expressed in terms of a single family of distributions.

2.1.3. Exponential family form and natural parameterisation

Suppose that $\mathbf{X} \sim \text{G-Wishart}(G, \delta, \mathbf{\Lambda})$. Then for \mathbf{X} such that $p(\mathbf{X}) > 0$ we have

$$p(\mathbf{X}) \propto \exp \left\{ \begin{bmatrix} \log|\mathbf{X}| \\ \text{vech}(\mathbf{X}) \end{bmatrix}^\top \begin{bmatrix} \frac{1}{2}(\delta - 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec}(\mathbf{\Lambda}) \end{bmatrix} \right\} = \exp\{\mathbf{T}(\mathbf{X})^\top \boldsymbol{\eta}\}, \tag{4}$$

where

$$\mathbf{T}(\mathbf{X}) = \begin{bmatrix} \log|\mathbf{X}| \\ \text{vech}(\mathbf{X}) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta} = \begin{bmatrix} \eta_1 \\ \boldsymbol{\eta}_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2}(\delta - 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec}(\mathbf{\Lambda}) \end{bmatrix}$$

are respectively sufficient statistic and natural parameter vectors. The inverse of the natural parameter mapping is

$$\begin{cases} \delta = 2(\eta_1 + 1), \\ \mathbf{\Lambda} = -2 \text{vec}^{-1}(\mathbf{D}_d^{+\top} \boldsymbol{\eta}_2). \end{cases}$$

Note that, throughout this article, we use $\text{vech}(\mathbf{X})$ rather than $\text{vec}(\mathbf{X})$ since the former is more compact and avoids duplications. Section S.1 in the web-supplement has further discussion on this matter.

2.2. The Inverse G-Wishart distribution

Suppose that $\mathbf{X} \sim \text{G-Wishart}(G, \delta, \mathbf{\Lambda})$, where \mathbf{X} is $d \times d$, and $\mathbf{Y} = \mathbf{X}^{-1}$. Let the density functions of \mathbf{X} and \mathbf{Y} be denoted by p_X and p_Y respectively. Then the density function of \mathbf{Y} is

$$p_Y(\mathbf{Y}) = p_X(\mathbf{Y}^{-1})|J(\mathbf{Y})|, \tag{5}$$

where

$$J(\mathbf{Y}) \equiv \text{the determinant of } \frac{\partial \text{vec}(\mathbf{Y}^{-1})}{\partial \text{vec}(\mathbf{Y})^\top}$$

is the Jacobian of the transformation.

An important observation is that the form of $J(\mathbf{Y})$ is dependent on the graph G . In the case of G being a decomposable graph an expression for $J(\mathbf{Y})$ is given by (2.4) of Letac & Massam (2007), with credit given to Roverato (2000). Therefore, if G is decomposable, the density function of an Inverse G -Wishart random matrix can be obtained by substitution of (2.4) of Letac & Massam (2007) into (5). However, depending on the complexity of G , simplification of the density function expression may be challenging.

With variational message passing in mind, we now turn to the $G = G_{\text{full}}$ and $G = G_{\text{diag}}$ special cases. The $G = G_{\text{diag}}$ case is simple since it involves products of univariate density functions and we have

$$\text{if } G = G_{\text{diag}} \text{ then } |J(\mathbf{Y})| = |\mathbf{Y}|^{-2}, \text{ for any } d \in \mathbb{N}. \tag{6}$$

The $G = G_{\text{full}}$ case is more challenging and is the focus of Muirhead (1982 Theorem 2.1.8):

$$\text{if } G = G_{\text{full}}, \text{ then } |J(\mathbf{Y})| = |\mathbf{Y}|^{-(d+1)}. \tag{7}$$

This result is also stated as Lemma 2.1 in Letac & Massam (2007).

Combining (5), (6) and (7) we have:

Result 3 Suppose that $\mathbf{Y} = \mathbf{X}^{-1}$, where $\mathbf{X} \sim G\text{-Wishart}(G, \delta, \mathbf{\Lambda})$ and \mathbf{X} is $d \times d$.

- (i) If $G = G_{\text{full}}$ then $p(\mathbf{Y}) \propto |\mathbf{Y}|^{-(\delta+2d)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{Y}^{-1}) \right\}$.
- (ii) If $G = G_{\text{diag}}$ then $p(\mathbf{Y}) \propto |\mathbf{Y}|^{-(\delta+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{Y}^{-1}) \right\}$.

While Result 3 only covers $G = G_{\text{full}}$ or $G = G_{\text{diag}}$ it shows that, in these special cases, the density function of an Inverse G -Wishart random matrix \mathbf{Y} is proportional to a power of $|\mathbf{Y}|$ multiplied by an exponentiated trace of a matrix multiplied by \mathbf{Y}^{-1} . This form does not necessarily arise for $G \notin \{G_{\text{full}}, G_{\text{diag}}\}$. Since the motivating variational message passing fragment update algorithms only involve the $G \in \{G_{\text{full}}, G_{\text{diag}}\}$ cases we focus on them for the remainder of this section.

2.2.1. The Inverse G -Wishart distribution when $G \in \{G_{\text{full}}, G_{\text{diag}}\}$

For succinct statement of variational message passing fragment update algorithms involving variance and covariance matrix parameters it is advantageous to have a single Inverse G -Wishart distribution notation for the $G \in \{G_{\text{full}}, G_{\text{diag}}\}$ cases.

Definition 2 Let \mathbf{X} be a $d \times d$ symmetric and positive definite random matrix and G be a d -node undirected graph such that \mathbf{X}^{-1} respects G . Let $\xi > 0$ and $\mathbf{\Lambda}$ be a symmetric positive definite $d \times d$ matrix $\mathbf{\Lambda}$.

- (i) If $G = G_{\text{full}}$ and ξ is restricted such that $\xi > 2d - 2$, then we say that \mathbf{X} has an Inverse G -Wishart distribution with graph G , shape parameter ξ and scale matrix $\mathbf{\Lambda}$, and write

$$\mathbf{X} \sim \text{Inverse } G\text{-Wishart}(G, \xi, \mathbf{\Lambda}),$$

if and only if the non-zero values of the density function of \mathbf{X} satisfy

$$p(\mathbf{X}) \propto |\mathbf{X}|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{X}^{-1}) \right\}.$$

- (ii) If $G = G_{\text{diag}}$, then say that \mathbf{X} has an Inverse G-Wishart distribution with graph G , shape parameter ξ and scale matrix $\mathbf{\Lambda}$, and write

$$\mathbf{X} \sim \text{Inverse G-Wishart}(G, \xi, \mathbf{\Lambda}),$$

if and only if the non-zero values of the density function of \mathbf{X} satisfy

$$p(\mathbf{X}) \propto |\mathbf{X}|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{X}^{-1}) \right\}.$$

- (iii) If $G \notin \{G_{\text{full}}, G_{\text{diag}}\}$, then $\mathbf{X} \sim \text{Inverse G-Wishart}(G, \xi, \mathbf{\Lambda})$ is not defined.

The shape parameter ξ used in Definition 2 is a reasonable compromise between various competing parameterisation choices for the Inverse G-Wishart distribution for $G \in \{G_{\text{full}}, G_{\text{diag}}\}$ and for use in variational message passing algorithms. It has the following attractions:

- The exponent of the determinant in the density function expression is $-(\xi + 2)/2$ regardless of whether $G = G_{\text{full}}$ or $G = G_{\text{diag}}$, which is consistent with the G-Wishart distributional notation used in Definition 1.
- In the $d = 1$ case ξ matches the shape parameter in the most common parameterisation of the Inverse Chi-Squared distribution such as that used in Gelman *et al.* (2014, Table A.1).

In the case where $\mathbf{X} \sim \text{Inverse G-Wishart}(G_{\text{full}}, \xi, \mathbf{\Lambda})$ we have the following:

Result 4 If the $d \times d$ random matrix \mathbf{X} is such that $\mathbf{X} \sim \text{Inverse G-Wishart}(G_{\text{full}}, \xi, \mathbf{\Lambda})$, then

$$p(\mathbf{X}) = \frac{|\mathbf{\Lambda}|^{(\xi-d+1)/2}}{2^{d(\xi-d+1)/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma\left(\frac{\xi-d-j}{2} + 1\right)} |\mathbf{X}|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{X}^{-1}) \right\} \\ \times \mathbf{1}(\mathbf{X} \text{ a symmetric and positive definite } d \times d \text{ matrix}).$$

The mean of \mathbf{X}^{-1} is

$$E(\mathbf{X}^{-1}) = (\xi - d + 1)\mathbf{\Lambda}^{-1}.$$

Result 4 follows directly from the fact that $\mathbf{X} \sim \text{Inverse G-Wishart}(G_{\text{full}}, \xi, \mathbf{\Lambda})$ if and only if \mathbf{X} has an Inverse Wishart distribution and established results for the density function and mean of this distribution given in, for example, Gelman *et al.* (2014, Table A.1).

We now deal with the $G = G_{\text{diag}}$ case.

Definition 3 Let x be a random variable. For $\delta > 0$ and $\lambda > 0$ we say that the random variable x has an Inverse Chi-Squared distribution with shape parameter δ and rate parameter λ , and write

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

if and only if $1/x \sim \chi^2(\delta, \lambda)$. If $x \sim \text{Inverse-}\chi^2(\delta, \lambda)$, then the density function of x is

$$p(x) = \frac{(\lambda/2)^{\delta/2}}{\Gamma(\delta/2)} x^{-(\delta+2)/2} \exp \left\{ -(\lambda/2)/x \right\} \mathbf{1}(x > 0).$$

Result 5 Suppose that the $d \times d$ random matrix \mathbf{X} is such that $\mathbf{X} \sim \text{Inverse-G-Wishart}(G_{\text{diag}}, \xi, \mathbf{\Lambda})$. Then the non-zero entries of \mathbf{X} satisfy

$$X_{jj} \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(\xi, \Lambda_{jj}), \quad 1 \leq j \leq d,$$

where Λ_{jj} is the j th diagonal entry of $\mathbf{\Lambda}$. The density function of \mathbf{X} is

$$\begin{aligned} p(\mathbf{X}) &= \frac{|\mathbf{\Lambda}|^{\xi/2}}{2^{d\xi/2}\Gamma(\xi/2)^d} |\mathbf{X}|^{-(\xi+2)/2} \exp\left\{-\frac{1}{2} \text{tr}(\mathbf{\Lambda}\mathbf{X}^{-1})\right\} \prod_{j=1}^d \mathbf{1}(X_{jj} > 0) \\ &= \frac{\prod_{j=1}^d \Lambda_{jj}^{\xi/2}}{2^{d\xi/2}\Gamma(\xi/2)^d} \prod_{j=1}^d X_{jj}^{-(\xi+2)/2} \exp\left\{-\frac{1}{2} \sum_{j=1}^d (\Lambda_{jj}/X_{jj})\right\} \prod_{j=1}^d \mathbf{1}(X_{jj} > 0). \end{aligned}$$

The mean of \mathbf{X}^{-1} is

$$E(\mathbf{X}^{-1}) = \xi\mathbf{\Lambda}^{-1} = \xi \text{diag}(1/\Lambda_{11}, \dots, 1/\Lambda_{dd}).$$

2.2.2. Natural parameter forms and sufficient statistic expectations

Suppose that $\mathbf{X} \sim \text{Inverse-G-Wishart}(G, \xi, \mathbf{\Lambda})$ where $G \in \{G_{\text{full}}, G_{\text{diag}}\}$. Then for \mathbf{X} such that $p(\mathbf{X}) > 0$,

$$p(\mathbf{X}) \propto \exp\left\{\left[\begin{array}{c} \log|\mathbf{X}| \\ \text{vech}(\mathbf{X}^{-1}) \end{array}\right]^{\top} \left[\begin{array}{c} -(\xi+2)/2 \\ -\frac{1}{2}\mathbf{D}_d^{\top} \text{vec}(\mathbf{\Lambda}) \end{array}\right]\right\} = \exp\{\mathbf{T}(\mathbf{X})^{\top} \boldsymbol{\eta}\},$$

where

$$\mathbf{T}(\mathbf{X}) = \left[\begin{array}{c} \log|\mathbf{X}| \\ \text{vech}(\mathbf{X}^{-1}) \end{array}\right] \quad \text{and} \quad \boldsymbol{\eta} = \left[\begin{array}{c} \eta_1 \\ \boldsymbol{\eta}_2 \end{array}\right] = \left[\begin{array}{c} -\frac{1}{2}(\xi+2) \\ -\frac{1}{2}\mathbf{D}_d^{\top} \text{vec}(\mathbf{\Lambda}) \end{array}\right] \tag{8}$$

are respectively sufficient statistic and natural parameter vectors. The inverse of the natural parameter mapping is

$$\begin{cases} \xi = -2\eta_1 - 2, \\ \mathbf{\Lambda} = -2 \text{vec}^{-1}(\mathbf{D}_d^{+\top} \boldsymbol{\eta}_2). \end{cases} \tag{9}$$

As explained in Section S.1 of the web-supplement, alternatives to (8) are those that use $\text{vec}(\mathbf{X})$ instead of $\text{vech}(\mathbf{X})$. Throughout this article, we use the more compact ‘vech’ form.

The following result is fundamental to succinct formulation of updates of covariance and variance parameter fragment updates for variational message passing:

Result 6. If \mathbf{X} is a $d \times d$ random matrix that has an Inverse G-Wishart distribution with graph $G \in \{G_{\text{full}}, G_{\text{diag}}\}$ and natural parameter vector $\boldsymbol{\eta}$. Then

$$E(\mathbf{X}^{-1}) = \begin{cases} \{\eta_1 + \frac{1}{2}(d+1)\}\{\text{vec}^{-1}(\mathbf{D}_d^{+\top} \boldsymbol{\eta}_2)\}^{-1}, & \text{if } G = G_{\text{full}}, \\ (\eta_1 + 1)\{\text{vec}^{-1}(\mathbf{D}_d^{+\top} \boldsymbol{\eta}_2)\}^{-1}, & \text{if } G = G_{\text{diag}}. \end{cases}$$

2.2.3. Relationships with the Hyper Inverse Wishart distributions

Throughout this article we follow the G-Wishart nomenclature as used by, for example, Atay-Kayis & Massam (2005), Letac & Massam (2007) and Uhler *et al.* (2018) in our naming of the Inverse G-Wishart family. Some earlier articles, such as Roverato (2000), use the term *Hyper Inverse Wishart* for the same family of distributions. The naming used here is in keeping with the more recent literature concerning Wishart distributions with graphical restrictions.

3. Connections with some recent covariance matrix distributions

Recently, Huang & Wand (2013) and Mulder & Pericchi (2018) developed covariance matrix distributional families that have attractions in terms of the types of marginal prior distributions that can be imposed on interpretable parameters within the covariance matrix. Mulder & Pericchi (2018) referred to their proposal as the *Matrix-F* family of distributions.

3.1. The Huang–Wand family of distributions

A major motivation for working with the Inverse G-Wishart distribution is the fact that the family of marginally non-informative priors proposed in Huang & Wand (2013) can be expressed succinctly in terms of the Inverse-G-Wishart(G, ζ, Λ) family where $G \in \{G_{\text{full}}, G_{\text{diag}}\}$. This means that variational message fragments that cater for Huang–Wand prior specification, as well as Inverse-Wishart prior specification, only require natural parameter vector manipulations within a single distributional family.

If Σ is a $d \times d$ symmetric positive definite matrix then, for $v_{\text{HW}} > 0$ and $s_1, \dots, s_d > 0$, the specification

$$\begin{aligned} \Sigma | \mathcal{A} &\sim \text{Inverse-G-Wishart}\left(G_{\text{full}}, v_{\text{HW}} + 2d - 2, \mathcal{A}^{-1}\right), \\ \mathcal{A} &\sim \text{Inverse-G-Wishart}\left(G_{\text{diag}}, 1, \{v_{\text{HW}} \text{diag}(s_1^2, \dots, s_d^2)\}^{-1}\right) \end{aligned} \tag{10}$$

places a distribution of the type given in Huang & Wand (2013) on Σ with shape parameter v_{HW} and scale parameters s_1, \dots, s_d .

The specification (10) matches (2) of Huang & Wand (2013) but with some differences in notation. First, d is used for matrix dimension here rather than p in Huang & Wand (2013). Also, the s_j , $1 \leq j \leq d$, scale parameters are denoted by A_j in Huang & Wand (2013). The a_j auxiliary variables in (2) of Huang & Wand (2013) are related to the matrix \mathcal{A} via the expression $\text{diag}(a_1, \dots, a_d) = 2v_{\text{HW}}\mathcal{A}$.

As discussed in Huang & Wand (2013), special cases of (10) correspond to marginally noninformative prior specification of the covariance matrix Σ in the sense that the standard deviation parameters $\sigma_j = (\Sigma)_{jj}^{1/2}$, $1 \leq j \leq d$, can have Half- t priors with arbitrarily large-scale parameters, controlled by the s_j values. This is in keeping with the advice given in Gelman (2006). Moreover, correlation parameters $\rho_{jj'} = (\Sigma)_{jj'}^{1/2} / (\sigma_j \sigma_{j'})$, for each $j \neq j'$ pair, have a Uniform distribution over the interval $(-1, 1)$ when $v_{\text{HW}} = 2$. We refer to this special case as the *Huang–Wand* marginally non-informative prior distribution with scale parameters s_1, \dots, s_d and write

$$\Sigma \sim \text{Huang–Wand}(s_1, \dots, s_d) \tag{11}$$

as a shorthand for (10) with $v_{\text{HW}} = 2$.

3.2. The Matrix- F family of distributions

For $v_{MP} > d - 1$, $\delta_{MP} > 0$ and \mathbf{B}_{MP} a $d \times d$ symmetric positive definite matrix, Mulder & Pericchi (2018) defined a $d \times d$ random matrix Σ to have a Matrix- F distribution, written

$$\Sigma \sim F(v_{MP}, \delta_{MP}, \mathbf{B}_{MP}), \tag{12}$$

if its density function has the form

$$p(\Sigma) \propto \frac{|\Sigma|^{(v_{MP}-d-1)/2} \mathbf{1}(\Sigma \text{ symmetric and positive definite})}{|\mathbf{I}_d + \Sigma \mathbf{B}_{MP}^{-1}|^{(v_{MP}+\delta_{MP}+d-1)/2}}.$$

However, standard manipulations of results given in Mulder & Pericchi (2018) show that specification (12) is equivalent to

$$\begin{aligned} \Sigma | \mathbf{A} &\sim \text{Inverse-G-Wishart}(G_{\text{full}}, \delta_{MP} + 2d - 2, \mathbf{A}^{-1}), \\ \mathbf{A} &\sim \text{Inverse-G-Wishart}(G_{\text{full}}, v_{MP} + d - 1, \mathbf{B}_{MP}^{-1}) \end{aligned} \tag{13}$$

in the notation used in the current paper. An important difference between (10) and (13) is that the former involves \mathbf{A} having an Inverse G-Wishart distribution with the restriction $G = G_{\text{diag}}$, while the latter has $G = G_{\text{full}}$. Section 2.4 of Mulder & Pericchi (2018) compares the two specifications in terms of the types of prior distributions that can be imposed on standard deviation and correlation parameters.

4. Variational message passing background

The overarching goal of this article was to identify and specify algebraic primitives for flexible imposition of covariance matrix priors within a variational message passing framework. In Wand (2017), these algebraic primitives are organised into fragments. This formalism is also used in Nolan & Wand (2017), Maestrini & Wand (2018) and McLean & Wand (2019).

Despite it being a central theme of this article, we will not provide a detailed description of the variational message passing here. Instead, we refer the reader to Wand (2017 Sections 2–4) for the relevant variational message passing background material.

Since the notational conventions for messages used in this section’s references are used in the remainder of this article, we summarize them here. If f denotes a generic factor and θ denotes a generic stochastic variable that is a neighbour of f in the factor graph then the message passed from f to θ and the message passed from θ to f are both functions of θ and are denoted by respectively,

$$m_{f \rightarrow \theta}(\theta) \quad \text{and} \quad m_{\theta \rightarrow f}(\theta).$$

Typically, the messages are proportional to an exponential family density function with sufficient statistic $\mathbf{T}(\theta)$, and we have

$$m_{f \rightarrow \theta}(\theta) \propto \exp \left\{ \mathbf{T}(\theta)^\top \boldsymbol{\eta}_{f \rightarrow \theta} \right\} \quad \text{and} \quad m_{\theta \rightarrow f}(\theta) \propto \exp \left\{ \mathbf{T}(\theta)^\top \boldsymbol{\eta}_{\theta \rightarrow f} \right\},$$

where $\boldsymbol{\eta}_{f \rightarrow \theta}$ and $\boldsymbol{\eta}_{\theta \rightarrow f}$ are the message natural parameter vectors. Such vectors play a central role in variational message passing iterative algorithms. We also adopt the notation

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



Figure 2. Diagram of the Inverse G-Wishart prior fragment.

5. The Inverse G-Wishart prior fragment

The Inverse G-Wishart prior fragment corresponds to the following prior imposition on a $d \times d$ covariance matrix Θ :

$$\Theta \sim \text{Inverse-G-Wishart}(G_\Theta, \xi_\Theta, \Lambda_\Theta),$$

for a d -node undirected graph G_Θ , scalar shape parameter ξ_Θ and scale matrix Λ_Θ . The fragment's factor is

$$p(\Theta) \propto |\Theta|^{-(\xi_\Theta+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Lambda_\Theta \Theta^{-1}) \right\} \\ \times \mathbf{1}(\Theta \text{ is symmetric and positive definite and } \Theta^{-1} \text{ respects } G_\Theta).$$

Figure 2 is a diagram of the fragment, which shows that its only factor to stochastic node message is

$$m_{p(\Theta) \rightarrow \Theta}(\Theta) \propto p(\Theta),$$

which leads to

$$m_{p(\Theta) \rightarrow \Theta}(\Theta) = \exp \left\{ \left[\begin{array}{c} \log |\Theta| \\ \text{vech}(\Theta^{-1}) \end{array} \right]^\top \left[\begin{array}{c} -\frac{1}{2}(\xi_\Theta + 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec}(\Lambda_\Theta) \end{array} \right] \right\}.$$

Therefore, the natural parameter update is

$$\eta_{p(\Theta) \rightarrow \Theta} \leftarrow \left[\begin{array}{c} -\frac{1}{2}(\xi_\Theta + 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec}(\Lambda_\Theta) \end{array} \right].$$

Apart from passing the natural parameter vector out of the fragment, we should also pass the graph out of the fragment. This entails the update:

$$G_{p(\Theta) \rightarrow \Theta} \leftarrow G_\Theta.$$

Algorithm 1 provides the inputs, updates and outputs for the Inverse G-Wishart prior fragment.

Algorithm 1 *The inputs, updates and outputs for the Inverse G-Wishart prior fragment.*

Hyperparameter Inputs: $G_\Theta, \xi_\Theta, \Lambda_\Theta$.
Updates:

$$\eta_{p(\Theta) \rightarrow \Theta} \leftarrow \left[\begin{array}{c} -\frac{1}{2}(\xi_\Theta + 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec}(\Lambda_\Theta) \end{array} \right]; G_{p(\Theta) \rightarrow \Theta} \leftarrow G_\Theta.$$

Outputs: $G_{p(\Theta) \rightarrow \Theta}, \eta_{p(\Theta) \rightarrow \Theta}$.

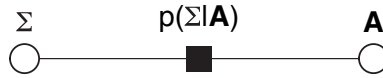


Figure 3. Diagram of the iterated Inverse G-Wishart fragment.

6. The Iterated Inverse G-Wishart fragment

The iterated Inverse G-Wishart fragment corresponds to the following specification involving a $d \times d$ covariance matrix Σ :

$$\Sigma | A \sim \text{Inverse-G-Wishart}(G, \xi, A^{-1}),$$

where G is a d -node undirected graph such that $G \in \{G_{\text{full}}, G_{\text{diag}}\}$ and ξ is a particular deterministic value of the Inverse G-Wishart shape parameter according to Definition 2. Figure 3 is a diagram of this fragment, showing that it has a factor $p(\Sigma | A)$ connected to two stochastic nodes Σ and A .

The factor of the iterated Inverse G-Wishart fragment is, as a function of both Σ and A ,

$$p(\Sigma | A) \propto \begin{cases} |A|^{-(\xi-d+1)/2} |\Sigma|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(A^{-1} \Sigma^{-1}) \right\}, & \text{if } G = G_{\text{full}}, \\ |A|^{-\xi/2} |\Sigma|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(A^{-1} \Sigma^{-1}) \right\}, & \text{if } G = G_{\text{diag}}. \end{cases}$$

As shown in Section S.2.1 of the web-supplement both of the factor to stochastic node messages of this fragment,

$$m_{p(\Sigma | A) \rightarrow \Sigma}(\Sigma) \quad \text{and} \quad m_{p(\Sigma | A) \rightarrow A}(A),$$

are proportional to Inverse G-Wishart density functions with graph $G \in \{G_{\text{full}}, G_{\text{diag}}\}$. We assume the following conjugacy constraints:

All messages passed to Σ and A from outside the fragment are proportional to Inverse G-Wishart density functions with graph $G \in \{G_{\text{full}}, G_{\text{diag}}\}$. The Inverse G-Wishart messages passed between Σ and $p(\Sigma | A)$ have the same graph. The Inverse G-Wishart messages passed between A and $p(\Sigma | A)$ have the same graph.

Under these constraints, and in view of for example, (7) of Wand (2017), the message passed from Σ to $p(\Sigma | A)$ has the form

$$m_{\Sigma \rightarrow p(\Sigma | A)}(\Sigma) = \exp \left\{ \left[\begin{array}{c} \log |\Sigma| \\ \text{vech}(\Sigma^{-1}) \end{array} \right]^\top \eta_{\Sigma \rightarrow p(\Sigma | A)} \right\}$$

and the message passed from A to $p(\Sigma | A)$ has the form

$$m_{A \rightarrow p(\Sigma | A)}(A) = \exp \left\{ \left[\begin{array}{c} \log |A| \\ \text{vech}(A^{-1}) \end{array} \right]^\top \eta_{A \rightarrow p(\Sigma | A)} \right\}.$$

Algorithm 2 gives the full set of updates of the message natural parameter vectors and graphs for the iterated Inverse-G-Wishart fragment. The derivation of Algorithm 2 is given in Section S.2 of the web-supplement.

Algorithm 2 *The inputs, updates and outputs for the iterated Inverse G-Wishart fragment.*

Graph Input: $G \in \{G_{\text{full}}, G_{\text{diag}}\}$.

Shape Parameter Input: $\xi > 0$.

Message Graph Input: $G_{A \rightarrow p(\Sigma|A)} \in \{G_{\text{full}}, G_{\text{diag}}\}$.

Natural Parameter Inputs: $\eta_{\Sigma \rightarrow p(\Sigma|A)}, \eta_{p(\Sigma|A) \rightarrow \Sigma}, \eta_{A \rightarrow p(\Sigma|A)}, \eta_{p(\Sigma|A) \rightarrow A}$.

Updates:

$$\begin{aligned}
 G_{p(\Sigma|A) \rightarrow \Sigma} &\leftarrow G; G_{p(\Sigma|A) \rightarrow A} \leftarrow G_{A \rightarrow p(\Sigma|A)} \\
 \eta_{p(\Sigma|A) \leftrightarrow \Sigma} &\leftarrow \eta_{p(\Sigma|A) \rightarrow \Sigma} + \eta_{\Sigma \rightarrow p(\Sigma|A)} \\
 \eta_{p(\Sigma|A) \leftrightarrow A} &\leftarrow \eta_{p(\Sigma|A) \rightarrow A} + \eta_{A \rightarrow p(\Sigma|A)} \\
 \text{If } G_{p(\Sigma|A) \rightarrow A} = G_{\text{full}}, &\text{ then } \omega_1 \leftarrow (d + 1)/2 \\
 \text{If } G_{p(\Sigma|A) \rightarrow A} = G_{\text{diag}}, &\text{ then } \omega_1 \leftarrow 1 \\
 E_q(\mathcal{A}^{-1}) &\leftarrow \left\{ \left(\eta_{p(\Sigma|A) \leftrightarrow A} \right)_1 + \omega_1 \right\} \left\{ \text{vec}^{-1} \left(\mathbf{D}_d^{+\top} \left(\eta_{p(\Sigma|A) \leftrightarrow A} \right)_2 \right) \right\}^{-1} \\
 \text{If } G_{p(\Sigma|A) \rightarrow \Sigma} = G_{\text{diag}}, &\text{ then } E_q(\mathcal{A}^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left(E_q(\mathcal{A}^{-1}) \right) \right\} \\
 \eta_{p(\Sigma|A) \rightarrow \Sigma} &\leftarrow \begin{bmatrix} -\frac{1}{2}(\xi + 2) \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec} \left(E_q(\mathcal{A}^{-1}) \right) \end{bmatrix} \\
 \text{If } G_{p(\Sigma|A) \rightarrow \Sigma} = G_{\text{full}}, &\text{ then } \omega_2 \leftarrow (d + 1)/2 \\
 \text{If } G_{p(\Sigma|A) \rightarrow \Sigma} = G_{\text{diag}}, &\text{ then } \omega_2 \leftarrow 1 \\
 E_q(\Sigma^{-1}) &\leftarrow \left\{ \left(\eta_{p(\Sigma|A) \leftrightarrow \Sigma} \right)_1 + \omega_2 \right\} \left\{ \text{vec}^{-1} \left(\mathbf{D}_d^{+\top} \left(\eta_{p(\Sigma|A) \leftrightarrow \Sigma} \right)_2 \right) \right\}^{-1} \\
 \text{If } G_{p(\Sigma|A) \rightarrow A} = G_{\text{diag}}, &\text{ then } E_q(\Sigma^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left(E_q(\Sigma^{-1}) \right) \right\} \\
 \eta_{p(\Sigma|A) \rightarrow A} &\leftarrow \begin{bmatrix} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2} \mathbf{D}_d^\top \text{vec} \left(E_q(\Sigma^{-1}) \right) \end{bmatrix}
 \end{aligned}$$

Outputs: $G_{p(\Sigma|A) \rightarrow \Sigma}, G_{p(\Sigma|A) \rightarrow A}, \eta_{p(\Sigma|A) \rightarrow \Sigma}, \eta_{p(\Sigma|A) \rightarrow A}$.

6.1. Corrections to Section 4.1.3 of Wand (2017)

The iterated Inverse G-Wishart fragment was introduced in Wand (2017) Section 4.1.3) and it is one of the five fundamental fragments of semiparametric regression given in Table 1. However, there are some errors due to the author of Wand (2017) failing to recognise particular subtleties regarding the Inverse G-Wishart distribution, as discussed in Section 2.2. We now point out misleading or erroneous aspects in Wand (2017) Section 4.1.3).

Firstly, in Wand (2017) Θ_1 plays the role of Σ and Θ_2 plays the role of A . The dimension of Θ_1 and Θ_2 is denoted by d^Θ . The first displayed equation of Wand (2017 Section 4.1.3) is

$$\Theta_1 | \Theta_2 \sim \text{Inverse-G-Wishart}(G, \kappa, \Theta_2^{-1}), \tag{14}$$

for $\kappa > d^\Theta - 1$ but it is only in the $G = G_{\text{full}}$ case that such a statement is reasonable for general $d^\Theta \in \mathbb{N}$. When $G = G_{\text{full}}$ then $\kappa = \xi - d^\Theta + 1$ according to the notation used in the current article. Therefore, (14) involves a different parameterisation to that used throughout this article. Therefore, our first correction is to replace the first displayed equation of Wand (2017 Section 4.1.3) by:

$$\Theta_1 | \Theta_2 \sim \text{Inverse-G-Wishart}(G, \xi, \Theta_2^{-1}),$$

Table 1. Specifications of inputs of Algorithms 1 and 2 for several variance, standard deviation and covariance matrix prior impositions.

Prior specification	Algorithm 1			Algorithm 2		
	G_{Θ}	ξ_{Θ}	Λ_{Θ}	ξ	G	$G_{\mathcal{A} \rightarrow p(\Sigma \mathcal{A})}$
$\sigma^2 \sim \text{Inverse-}\chi^2(\delta_{\sigma^2}, \lambda_{\sigma^2})$	G_{full}	δ_{σ^2}	λ_{σ^2}	N.A.	N.A.	N.A.
$\sigma^2 \sim \text{Inv:-Gamma}(\alpha_{\sigma^2}, \beta_{\sigma^2})$	G_{full}	$2\alpha_{\sigma^2}$	$2\beta_{\sigma^2}$	N.A.	N.A.	N.A.
$\Sigma \sim \text{Inv:-Wishart}(\kappa_{\Sigma}, \Lambda_{\Sigma})$	G_{full}	$\kappa_{\Sigma} + d - 1$	Λ_{Σ}	N.A.	N.A.	N.A.
$\sigma \sim \text{Half-}t(s_{\sigma}, \nu_{\sigma})$	G_{diag}	1	$(\nu_{\sigma} s_{\sigma}^2)^{-1}$	ν_{σ}	G_{full}	G_{diag}
$\sigma \sim \text{Half-Cauchy}(s_{\sigma})$	G_{diag}	1	$(s_{\sigma}^2)^{-1}$	1	G_{full}	G_{diag}
$\Sigma \sim \text{Huang-Wand}$ $(s_{\Sigma,1}, \dots, s_{\Sigma,d})$	G_{diag}	1	$\left\{ 2 \text{diag}(s_{\Sigma,1}^2, \dots, s_{\Sigma,d}^2) \right\}^{-1}$	$2d$	G_{full}	G_{diag}
$\Sigma \sim \text{Matrix-F}$ $(\nu_{\text{MP}}, \delta_{\text{MP}}, \mathbf{B}_{\text{MP}})$	G_{full}	$\nu_{\text{MP}} + d - 1$	$\mathbf{B}_{\text{MP}}^{-1}$	$\frac{\delta_{\text{MP}} + 2d - 2}{2d - 2}$	G_{full}	G_{full}

Note: The abbreviation N.A. stands for not applicable since Algorithm 2 is not needed for the first three prior impositions.

where $\xi > 0$ if $G = G_{\text{diag}}$ and $\xi > 2d^{\Theta} - 2$ if $G = G_{\text{full}}$.

The following sentence in Wand (2017) Section 4.1.3): ‘The fragment factor is of the form

$$p(\Theta_1 | \Theta_2) \propto |\Theta_2|^{-\kappa/2} |\Theta_1|^{-(\kappa+d^{\Theta}+1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Theta_1^{-1} \Theta_2^{-1}) \right\},$$

should instead be ‘The fragment factor is of the form

$$p(\Theta_1 | \Theta_2) \propto \begin{cases} |\Theta_2|^{-(\xi-d^{\Theta}+1)/2} |\Theta_1|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Theta_1^{-1} \Theta_2^{-1}) \right\} & \text{if } G = G_{\text{full}}, \\ |\Theta_2|^{-\xi/2} |\Theta_1|^{-(\xi+2)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Theta_1^{-1} \Theta_2^{-1}) \right\} & \text{if } G = G_{\text{diag}}. \end{cases}$$

In (31) of Wand (2017), the first entry of the vector on the right-hand side of the \leftarrow should be

$$-(\xi + 2)/2 \quad \text{rather than} \quad -(\kappa + d^{\Theta} + 1)/2.$$

To match the correct parameterisation of the Inverse G-Wishart distribution, as used in the current article, (32) of Wand (2017) should be

$$\text{‘E}(X^{-1}) \quad \text{where} \quad X \sim \text{Inverse-G-Wishart}(G, \xi, \Lambda)\text{’}.$$

The equation in Wand (2017) Section 4.1.3):

$$\text{‘}\eta_{p(\Theta_1|\Theta_2) \rightarrow \Theta_2} \leftarrow \left[\begin{array}{c} -\kappa/2 \\ -\frac{1}{2} \text{vec} \left(\text{E}_{p(\Theta_1|\Theta_2) \rightarrow \Theta_2}(\Theta_1^{-1}) \right) \end{array} \right]\text{’},$$

should be replaced by

$$\text{‘}\eta_{p(\Theta_1|\Theta_2) \rightarrow \Theta_2} \leftarrow \left[\begin{array}{c} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2} \text{vec} \left(\text{E}_{p(\Theta_1|\Theta_2) \rightarrow \Theta_2}(\Theta_1^{-1}) \right) \end{array} \right]\text{’},$$

where ω_2 depends on the graph of the Inverse G-Wishart distribution corresponding to $\text{E}_{p(\Theta_1|\Theta_2) \rightarrow \Theta_2}$. If the graph is G_{full} then $\omega_2 = (d^{\Theta} + 1)/2$ and if the graph is G_{diag} then $\omega_2 = 1$.

Lastly the iterated Inverse G-Wishart fragment natural parameter updates given by (36) and (37) of Wand (2017) are affected by the oversights described in the preceding paragraphs. They should be replaced by the updates given in Algorithm 2 with $\Theta_1 = \Sigma$ and $\Theta_2 = A$.

7. Use of the fragments for covariance matrix prior specification

The underlying rationale for the Inverse G-Wishart prior and iterated Inverse G-Wishart fragments is their ability to facilitate the specification of a wide range of covariance matrix priors within the variational message passing framework. In the $d = 1$ special case, covariance matrix parameters reduce to variance parameters and their square roots are standard deviation parameters. In this section, we spell out how the fragments, and their natural parameter updates in Algorithms 1 and 2, can be used for prior specification in important special cases.

7.1. Imposing an Inverse Chi-Squared prior on a variance parameter

Let σ^2 be a variance parameter and consider the prior imposition

$$\sigma^2 \sim \text{Inverse-}\chi^2(\delta_{\sigma^2}, \lambda_{\sigma^2}),$$

for hyperparameters $\delta_{\sigma^2}, \lambda_{\sigma^2} > 0$, within a variational message passing scheme. Then Algorithm 1 should be called with inputs set to:

$$G_{\Theta} = G_{\text{full}}, \quad \xi_{\Theta} = \delta_{\sigma^2}, \quad \Lambda_{\Theta} = \lambda_{\sigma^2}.$$

7.2. Imposing an Inverse Gamma prior on a variance parameter

Let σ^2 be a variance parameter and consider the prior imposition

$$\sigma^2 \sim \text{Inverse-Gamma}(\alpha_{\sigma^2}, \beta_{\sigma^2}), \tag{15}$$

for hyperparameters $\alpha_{\sigma^2}, \beta_{\sigma^2} > 0$. The density function corresponding to (15) is

$$p(\sigma^2; \alpha_{\sigma^2}, \beta_{\sigma^2}) \propto (\sigma^2)^{-\alpha_{\sigma^2}-1} \exp\{-\beta_{\sigma^2}/(\sigma^2)\} \mathbf{1}(\sigma^2 > 0).$$

Note that the Inverse Chi-Squared and Inverse Gamma distributions are simple reparameterisations of each other since

$$x \sim \text{Inverse-}\chi^2(\delta, \lambda) \quad \text{if and only if} \quad x \sim \text{Inverse-Gamma}\left(\frac{1}{2}\delta, \frac{1}{2}\lambda\right).$$

To achieve (15) Algorithm 1 should be called with inputs set to:

$$G_{\Theta} = G_{\text{full}}, \quad \xi_{\Theta} = 2\alpha_{\sigma^2}, \quad \Lambda_{\Theta} = 2\beta_{\sigma^2}.$$

7.3. Imposing an Inverse Wishart prior on a covariance matrix parameter

A random matrix X is defined to have an Inverse Wishart distribution with shape parameter κ and scale matrix Σ , written $X \sim \text{Inverse-Wishart}(\kappa, \Lambda)$, if and only if the density function of X is

$$p(\mathbf{X}) = \frac{|\Lambda|^{\kappa/2}}{2^{\kappa d/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma\left(\frac{\kappa+1-j}{2}\right)} |\mathbf{X}|^{-(\kappa+d+1)/2} \exp\left\{-\frac{1}{2} \text{tr}(\Lambda \mathbf{X}^{-1})\right\} \quad (16)$$

$\times \mathbf{1}(\mathbf{X} \text{ a symmetric and positive definite } d \times d \text{ matrix}).$

Note that this is the common parameterisation of the Inverse Wishart distribution (e.g. Gelman *et al.* 2014, Table A.1). Crucially, (16) uses a *different* shape parametrisation from that used for the Inverse G-Wishart distribution in Definition 2 when $G = G_{\text{full}}$ with the relationship between the two shape parameters given by $\kappa = \xi - d + 1$. Even though the more general Inverse G-Wishart family is important for the internal workings of variational message passing, the ordinary Inverse Wishart distribution, with the parameterisation as given in (16), is more common when imposing a prior on a covariance matrix.

Let Σ be a $d \times d$ matrix and consider the prior imposition

$$\Sigma \sim \text{Inverse-Wishart}(\kappa_{\Sigma}, \Lambda_{\Sigma}) \quad (17)$$

for hyperparameters $\kappa_{\Sigma}, \Lambda_{\Sigma} > 0$, within a variational message passing scheme. Then Algorithm 1 should be called with inputs set to:

$$G_{\Theta} = G_{\text{full}}, \quad \xi_{\Theta} = \kappa_{\Sigma} + d - 1, \quad \Lambda_{\Theta} = \Lambda_{\Sigma}.$$

7.4. Imposing a Half- t prior on a standard deviation parameter

Consider the prior imposition

$$\sigma \sim \text{Half-}t(s_{\sigma}, v_{\sigma}), \quad (18)$$

for a scale parameter $s_{\sigma} > 0$ and a degrees of freedom parameter $v_{\sigma} > 0$. The density function corresponding to (18) is such that $p(\sigma) \propto \{1 + (\sigma/s_{\sigma})^2/v_{\sigma}\}^{-(v_{\sigma}+1)/2} \mathbf{1}(\sigma > 0)$. This is equivalent to

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(v_{\sigma}, 1/a) \quad \text{and} \quad a \sim \text{Inverse-}\chi^2(1, 1/s_{\sigma}^2). \quad (19)$$

Since $d = 1$, the graphs G_{full} and G_{diag} are the same – a single node graph. Treating σ^2 and a as 1×1 matrices we can re-write (19) as

$$\sigma^2 | a \sim \text{Inverse-G-Wishart}(G_{\text{full}}, v_{\sigma}, a^{-1}) \quad \text{and} \quad a \sim \text{Inverse-G-Wishart}(G_{\text{diag}}, 1, (v_{\sigma} s_{\sigma}^2)^{-1})$$

(e.g. Armagan, Dunson & Clyde 2011). The specification

$$a \sim \text{Inverse-G-Wishart}(G_{\text{diag}}, 1, (v_{\sigma} s_{\sigma}^2)^{-1})$$

involves calling Algorithm 1 with

$$G_{\Theta} = G_{\text{diag}}, \quad \xi_{\Theta} = 1 \quad \text{and} \quad \Lambda_{\Theta} = (v_{\sigma} s_{\sigma}^2)^{-1}.$$

The output is the single node graph $G_{p(\Theta) \rightarrow \Theta}$ and the 2×1 natural parameter vector

$$\eta_{p(\Theta) \rightarrow \Theta} = \eta_{p(a) \rightarrow a}.$$

The specification

$$\sigma^2 | a \sim \text{Inverse-G-Wishart}(G_{\text{full}}, v_{\sigma}, a^{-1})$$

implies that Algorithm 2 is called with graph input $G = G_{\text{full}}$, shape parameter input $\xi = v_\sigma$ and message parameter inputs

$$\boldsymbol{\eta}_{p(\Sigma|A)\rightarrow\Sigma} = \boldsymbol{\eta}_{p(\sigma^2|a)\rightarrow\sigma^2}, \quad \boldsymbol{\eta}_{\Sigma\rightarrow p(\Sigma|A)} = \boldsymbol{\eta}_{\sigma^2\rightarrow p(\sigma^2|a)},$$

and

$$G_{p(\Sigma|A)\rightarrow A} = G_{\text{diag}}, \quad \boldsymbol{\eta}_{p(\Sigma|A)\rightarrow A} = \boldsymbol{\eta}_{p(\sigma^2|a)\rightarrow a} \quad \text{and} \quad \boldsymbol{\eta}_{A\rightarrow p(\Sigma|A)} = \boldsymbol{\eta}_{a\rightarrow p(\sigma^2|a)}.$$

Note that in this $d = 1$ special case G_{full} and G_{diag} are both the single node graph.

7.4.1. The Half-Cauchy special case

The special case of

$$\sigma \sim \text{Half-Cauchy}(s_\sigma) \tag{20}$$

corresponds to $v_\sigma = 1$. The density function corresponding to (20) is such that $p(\sigma) \propto \{1 + (\sigma/s_\sigma)^2\}^{-1} \mathbf{1}(\sigma > 0)$. Therefore, one should set $\xi = 1$ in the call to Algorithm 2.

7.5. Imposing a Huang–Wand prior on a covariance matrix

To impose the Huang–Wand prior

$$\Sigma \sim \text{Huang–Wand}(s_{\Sigma,1}, \dots, s_{\Sigma,d})$$

in a variational message passing framework we should have the inputs to Algorithm 1 being as follows:

$$G_\Theta = G_{\text{diag}}, \quad \xi_\Theta = 1 \quad \text{and} \quad \Lambda_\Theta = \{2 \text{diag}(s_{\Sigma,1}^2, \dots, s_{\Sigma,d}^2)\}^{-1}.$$

The graph parameter input to Algorithm 2 should be $G = G_{\text{full}}$ and the shape parameter input should be $\xi = 2d$.

7.6. Imposing a Matrix- F prior on a covariance matrix

To impose the Matrix- F prior

$$\Sigma \sim F(v_{\text{MP}}, \delta_{\text{MP}}, \mathbf{B}_{\text{MP}})$$

in a variational message passing framework the inputs to Algorithm 1 should be as follows:

$$G_\Theta = G_{\text{full}}, \quad \xi_\Theta = v_{\text{MP}} + d - 1 \quad \text{and} \quad \Lambda_\Theta = \mathbf{B}_{\text{MP}}^{-1}.$$

The graph parameter input to Algorithm 2 should be $G = G_{\text{full}}$ and the shape parameter input should be $\xi = \delta_{\text{MP}} + 2d - 2$.

7.7. Tabular summary of fragment-based prior specification

Table 1 summarizes the results of this section and is a crucial reference for placing priors of covariance matrix, variance and standard deviation parameters in variational message passing schemes that make use of Algorithms 1 and 2.

8. Illustrative example

We illustrate the use of Algorithms 1 and 2 for the case of Bayesian linear mixed models with t distribution responses. Such t -based models impose a form of robustness in situations where the responses are susceptible to having outlying values (e.g. Lange, Little & Taylor 1989). The notation $y \sim t(\mu, \sigma, \nu)$ indicates that the random variable y has a t distribution with location parameter μ , scale parameter $\sigma > 0$ and degrees of freedom parameter $\nu > 0$. The corresponding density function of y is

$$p(y) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sigma\sqrt{\pi\nu}\Gamma(\nu/2)[1 + \{(y - \mu)/\sigma\}^2/\nu]^{\frac{\nu+1}{2}}}.$$

Now suppose that the response data consists of repeated measures within each of m groups. Let

$$y_{ij} = \text{the } j\text{th response for the } i\text{th group, } 1 \leq j \leq n_i, \quad 1 \leq i \leq m,$$

and then let \mathbf{y}_i , $1 \leq i \leq m$, be the $n_i \times 1$ vectors containing y_{ij} data for the i th group. For each $1 \leq i \leq m$, let \mathbf{X}_i be $n_i \times p$ design matrices corresponding to the fixed effects and \mathbf{Z}_i be $n_i \times q$ design matrices corresponding to the random effects. Next put

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{bmatrix} \quad \text{and} \quad \mathbf{Z} = \text{blockdiag}(\mathbf{Z}_i), \quad \substack{1 \leq i \leq m} \quad (21)$$

and define $N = n_1 + \dots + n_m$ to be the number of rows in each of \mathbf{y} , \mathbf{X} and \mathbf{Z} . Let y_ℓ be the ℓ th entry of \mathbf{y} , $1 \leq \ell \leq N$. The family of Bayesian t response linear mixed models that we consider is

$$\begin{aligned} y_\ell | \boldsymbol{\beta}, \mathbf{u}, \sigma &\stackrel{\text{ind.}}{\sim} t((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_\ell, \sigma, \nu), \quad 1 \leq \ell \leq N, \quad \mathbf{u} | \boldsymbol{\Sigma} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}), \\ \boldsymbol{\beta} &\sim \mathbf{N}(0, \sigma_\beta^2 \mathbf{I}), \quad \sigma \sim \text{Half-Cauchy}(s_\sigma), \quad \frac{1}{2}\nu \sim \text{Moon-Rock}(0, \lambda_\nu), \\ \boldsymbol{\Sigma} &\sim \text{Huang-Wand}(s_{\boldsymbol{\Sigma},1}, \dots, s_{\boldsymbol{\Sigma},q}), \end{aligned} \quad (22)$$

for hyperparameters $\sigma_\beta, s_\sigma, \lambda_\nu, s_{\boldsymbol{\Sigma},1}, \dots, s_{\boldsymbol{\Sigma},q} > 0$.

As explained in McLean & Wand (2019), the Moon Rock family of distributions is conjugate for the parameter $\frac{1}{2}\nu$, with the notation $x \sim \text{Moon-Rock}(\alpha, \beta)$ indicating that the corresponding density function satisfies $p(x) \propto \{x^x/\Gamma(x)\}^\alpha \exp(-\beta x)\mathbf{1}(x > 0)$. In the variational message passing treatment of the degrees of freedom parameter it is simpler to work with

$$v = \frac{1}{2}\nu, \quad \text{so that} \quad v \sim \text{Moon-Rock}(0, \lambda_\nu).$$

After the approximate posterior density function of v is obtained via variational message passing, it is trivial to then obtain the same for ν . Hence, we work with v , rather than ν , in the upcoming description of variational message passing-based fitting and inference for (22).

Next note that

$$y_\ell | \boldsymbol{\beta}, \mathbf{u}, \sigma \stackrel{\text{ind.}}{\sim} t((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_\ell, \sigma, 2v), \quad 1 \leq \ell \leq N,$$

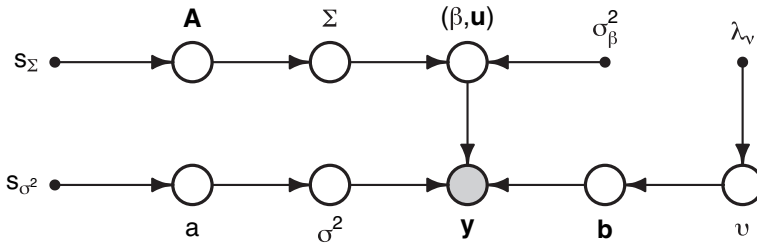


Figure 4. Directed acyclic graph corresponding to the t response linear mixed model (22) with auxiliary variable representations (23)–(25). The shaded circle corresponds to the observed data. The unshaded circles correspond to model parameters and auxiliary variables. The small solid circles correspond to hyperparameters.

is equivalent to

$$y_\ell | \boldsymbol{\beta}, \mathbf{u}, \sigma^2, b_\ell \sim N((X\boldsymbol{\beta} + \mathbf{Zu})_\ell, b_\ell \sigma^2), \quad b_\ell | v \overset{\text{ind.}}{\sim} \text{Inverse-}\chi^2(2v, 2v), \quad (23)$$

$\sigma \sim \text{Half-Cauchy}(s_\sigma)$ is equivalent to

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a) \quad \text{and} \quad a \sim \text{Inverse-}\chi^2(1, 1/s_\sigma^2) \quad (24)$$

and $\boldsymbol{\Sigma} \sim \text{Huang-Wand}(s_{\boldsymbol{\Sigma},1}, \dots, s_{\boldsymbol{\Sigma},q})$ is equivalent to

$$\begin{aligned} \boldsymbol{\Sigma} | \mathbf{A} &\sim \text{Inverse-G-Wishart}(G_{\text{full}}, 2q, \mathbf{A}^{-1}), \\ \mathbf{A} &\sim \text{Inverse-G-Wishart}(G_{\text{diag}}, 1, \{2 \text{diag}(s_{\boldsymbol{\Sigma},1}^2, \dots, s_{\boldsymbol{\Sigma},q}^2)\}^{-1}). \end{aligned} \quad (25)$$

Substitution of (23), (24) and (25) into (22) leads to the hierarchical Bayesian model depicted as a directed acyclic graph in Figure 4 with $\mathbf{b} = (b_1, \dots, b_N)$. The unshaded circles in Figure 4 correspond to model parameters and auxiliary variables and will be referred to as *hidden nodes*.

Consider the following mean-field approximation of the joint posterior of the hidden nodes in Figure 4

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b} | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}, a, \mathbf{A}, \mathbf{b}) q(\sigma^2, \boldsymbol{\Sigma}, v), \quad (26)$$

where q denotes the approximate posterior density functions of the relevant parameters. Application of induced factor results (e.g. Bishop 2006, Section 10.2.5) leads to the additional factorisations

$$q(\boldsymbol{\beta}, \mathbf{u}, a, \mathbf{A}, \mathbf{b}) = q(\boldsymbol{\beta}, \mathbf{u}) q(a) q(\mathbf{A}) \prod_{\ell=1}^N q(b_\ell) \quad \text{and} \quad q(\sigma^2, \boldsymbol{\Sigma}, v) = q(\sigma^2) q(\boldsymbol{\Sigma}) q(v),$$

and so the restriction given in (26) is equivalent to

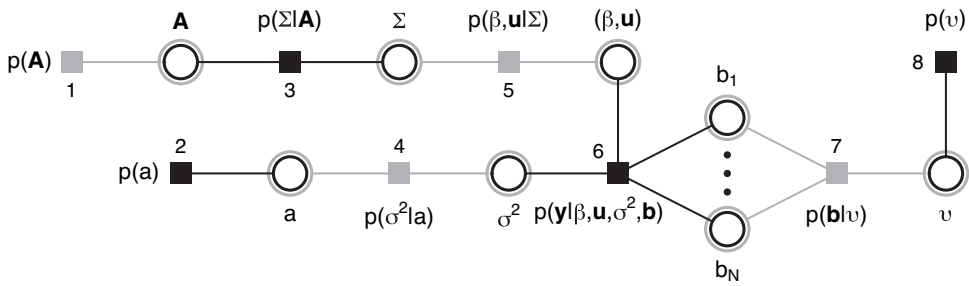


Figure 5. Factor graph corresponding to the t response linear mixed model (22) with auxiliary variable representations (23)–(25). The circular nodes correspond to stochastic nodes in the q -density factorisation in (27). The rectangular nodes correspond to the factors on the right-hand side of (28). The fragments are numbered 1 to 8 according to appearance from left to right. Shading is used to show the distinction between adjacent fragments.

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b}|\mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u})q(\sigma^2)q(v)q(\boldsymbol{\Sigma})q(a)q(\mathbf{A}) \prod_{\ell=1}^N q(b_\ell). \tag{27}$$

Figure 5 is a factor graph representation of the joint density function of all random variables and vectors, or stochastic nodes, in Figure 4 hierarchical model, with unshaded circles for each stochastic node according to the q -density factorisation given in (27) and filled-in rectangles corresponding to factors on the right-hand side of

$$p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b}) = p(\mathbf{A})p(a)p(\boldsymbol{\Sigma}|\mathbf{A})p(\sigma^2|a)p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma})p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{u}, \sigma^2, \mathbf{b})p(\mathbf{b}|v)p(v). \tag{28}$$

Edges join each factor to a stochastic node that appears in the factor. To aid upcoming discussion, the fragments are numbered 1 to 8 according to appearance from left to right. Recall that a fragment is a sub-graph consisting of a factor and all of its neighbouring nodes. Figure 5 uses shading to show the distinction between adjacent fragments.

Note that (e.g. Minka 2005; Wand 2017) the variational message passing iteration loop has the following generic steps:

1. Choose a factor.
2. Update the parameter vectors of the messages passed from the factor’s neighbouring stochastic nodes to the factor.
3. Update the parameter vectors of the messages passed from the factor to its neighbouring stochastic nodes.

Step 2. is very simple and has generic form given by, for example, (7) of Wand (2017). In the Figure 5 factor graph an example of Step 2. is:

$$\begin{aligned} &\text{the message passed from } \boldsymbol{\Sigma} \text{ to } p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma}) \\ &= \text{the message passed from } p(\boldsymbol{\Sigma}|\mathbf{A}) \text{ to } \boldsymbol{\Sigma} \text{ in the previous iteration.} \end{aligned} \tag{29}$$

In terms of natural parameter vector updates, (29) corresponds to:

$$\boldsymbol{\eta}_{\boldsymbol{\Sigma} \rightarrow p(\boldsymbol{\beta}, \mathbf{u}|\boldsymbol{\Sigma})} \longleftarrow \boldsymbol{\eta}_{p(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}.$$

Most of the other stochastic node to factor updates in Figure 5 have an analogous form. The exception are the messages passed within fragments 6 and 7, which require use of the slightly more complicated form as given by, for example, (7) of Wand (2017).

It remains to discuss Step 3, corresponding to the factor to stochastic node updates:

- Fragments 1 and 2 are Inverse G-Wishart prior fragments and the factor to stochastic node parameter vector updates are performed according to Algorithm 1. In view of Table 1, the graph and shape hyperparameter inputs are $G_{\Theta} = G_{\text{diag}}$ and $\xi_{\Theta} = 1$. For fragment 1 the rate hyperparameter is $\Lambda_{\Theta} = \{2 \text{diag}(s_{\Sigma,1}^2, \dots, s_{\Sigma,q}^2)\}^{-1}$. For fragment 2 the rate hyperparameter is $\Lambda_{\Theta} = (s_{\sigma}^2)^{-1}$.
- Fragments 3 and 4 are iterated Inverse G-Wishart prior fragments and the factor to stochastic node parameter vector updates are performed according to Algorithm 2. As shown in Table 1, the graph inputs should be

$$G_{A \rightarrow p(\Sigma|A)} = G_{\text{diag}}, \quad G_{a \rightarrow p(\sigma^2|a)} = G_{\text{diag}},$$

$$G_{\Sigma \rightarrow p(\Sigma|A)} = G_{\text{full}}, \quad \text{and} \quad G_{\sigma^2 \rightarrow p(\sigma^2|a)} = G_{\text{full}}.$$

The first two of these are imposed by the messages passed from fragments 1 and 2. For fragment 3, the shape parameter input is $\xi = 2q$. For fragment 4, the shape parameter input is $\xi = 1$.

- Fragment 5 is the Gaussian penalisation fragment described in Section 4.1.4 of Wand (2017) with, in the notation given there, $L = 1$, $\mu_{\theta_0} = \mathbf{0}$ and $\Sigma_{\theta_0} = \sigma_{\beta}^2 \mathbf{I}$.
- Fragments 6 and 7 correspond to the t likelihood fragment. Its natural parameter updates are provided by McLean & Wand (2019) Algorithm 2).
- Fragment 8 corresponds to the imposition of a Moon Rock prior distribution on a shape parameter. This is a very simple fragment for which the only inputs are the Moon Rock prior specification hyperparameters and the output is the natural parameter vector of the Moon Rock prior density function. Since this fragment is not listed as an algorithm in this article or elsewhere, we provide further details in the paragraph after the next one.

For Fragments 5, 6 and 7 simple conversions between two different versions of natural parameter vectors need to be made. Section S.1 of the web-supplement explains these conversions.

The most general Moon Rock prior specification for a generic parameter θ is

$$\theta \sim \text{Moon-Rock}(\alpha_{\theta}, \beta_{\theta}).$$

This corresponds to the prior density function having exponential family form

$$p(\theta) \propto \exp \left\{ \left[\begin{array}{c} \theta \log(\theta) - \log \Gamma(\theta) \\ \theta \end{array} \right]^{\top} \left[\begin{array}{c} \alpha_{\theta} \\ -\beta_{\theta} \end{array} \right] \right\}.$$

The inputs of the Moon Rock prior fragment are $\alpha_{\theta} \geq 0$ and $\beta_{\theta} > 0$ and the output is the natural parameter vector

$$\boldsymbol{\eta}_{p(\theta) \rightarrow \theta} \longleftarrow \left[\begin{array}{c} \alpha_{\theta} \\ -\beta_{\theta} \end{array} \right].$$

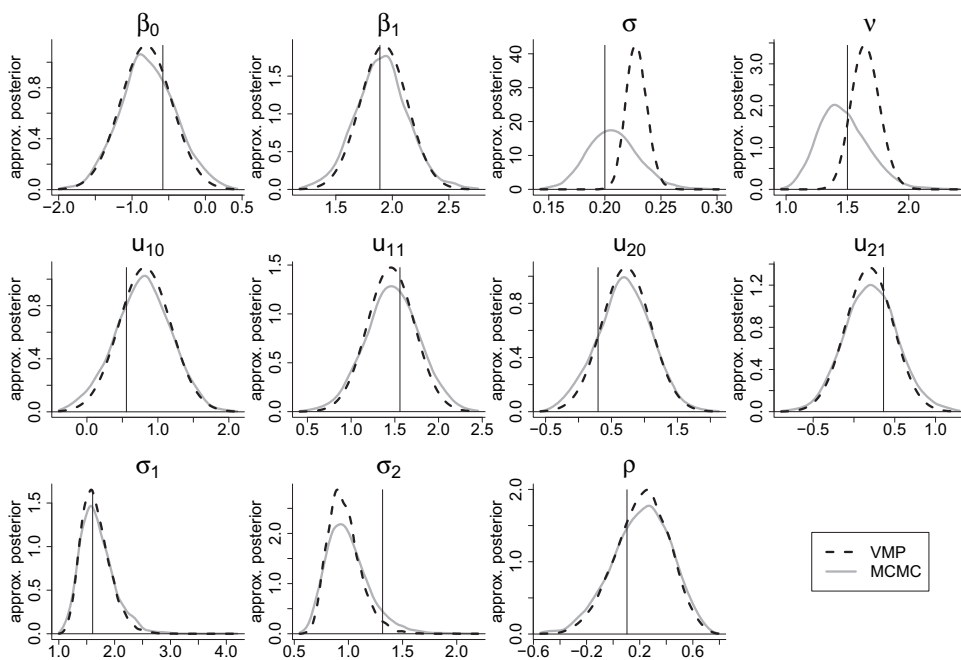


Figure 6. Approximate posterior density functions for the parameters in model (22) based on both variational message passing (VMP) and Markov chain Monte Carlo (MCMC) algorithms applied to data simulated according to the true values (30) and sample sizes, predictor values and hyperparameter values as described in the text. The vertical lines indicate true parameter values.

Since, for the t response mixed model illustrative example, we have the prior imposition $v \sim \text{Moon-Rock}(0, \lambda_v)$ we simply call the Moon Rock prior fragment with $(\alpha_\theta, \beta_\theta)$ set to $(0, \lambda_v)$.

To demonstrate variational message passing for fitting and inference for model (22), we simulated data according to the dimension values $p = q = 2$ and the true parameter values

$$\beta_{\text{true}} = \begin{bmatrix} -0.58 \\ 1.89 \end{bmatrix}, \quad \sigma_{\text{true}}^2 = 0.2, \quad \Sigma_{\text{true}} = \begin{bmatrix} 2.58 & 0.22 \\ 0.22 & 1.73 \end{bmatrix} \quad \text{and} \quad v_{\text{true}} = 1.5. \quad (30)$$

The sample sizes were $m = 20$, with $n_i = 15$ observations per group, and the predictor data were generated from the Uniform distribution on the unit interval. The hyperparameter values were set at

$$\sigma_\beta = s_\sigma = s_{\Sigma,1} = s_{\Sigma,2} = 10^5 \quad \text{and} \quad \lambda_v = 0.01.$$

We ran the variational message passing algorithm as described above until the relative change the variational parameters is below 10^{-10} . as well as Markov chain Monte Carlo via the R language (R Core Team 2021) package rstan (Stan Development Team 2020). For Markov chain Monte Carlo fitting, a warmup of size 1000 was used, followed by chains of size 5000 retained for inference.

Figure 6 compares the approximate posterior density functions based on both variational message passing (VMP) and Markov chain Monte Carlo (MCMC). The middle row performs

the comparison for the random intercept and slope parameters, u_{i0} and u_{i1} , for $i = 1, 2$. The parameters in the third row of Figure 6 are for the standard deviation and correlation parameters in the Σ matrix, according to the notation $(\Sigma)_{11} = \sigma_1$, $(\Sigma)_{22} = \sigma_2$ and $(\Sigma)_{12} = \sigma_1\sigma_2\rho$. For most of the stochastic nodes, the accuracy of variational message passing is seen to be very good. For σ and ν , some under-approximation of the spread and locational shift is apparent. A likely root cause is the imposition of the product restriction $q(\sigma, \nu) = q(\sigma)q(\nu)$ even though these two parameters have a significant amount of posterior dependence.

We have prepared a bundle of R language code that carries out variational message passing for this illustrative example, including use of Algorithms 1 and 2 for the imposition of Half Cauchy and Huang–Wand priors. This code is part of the web-supplement for this article.

Lastly, we point out that this illustrative example does not involve matrix algebraic streamlining for random effects models. This relatively new area for variational message passing research, which streamlines calculations involving sparse matrix forms that arise in linear mixed models, is described in Nolan, Menictas & Wand (2020).

9. Closing remarks

Algorithms 1 and, especially, Algorithm 2 and their underpinnings are quite involved and dependent upon a careful study of particular special cases of the inverses of G-Wishart random matrices. The amount of detail provided by this article is tedious, but necessary, to ensure that the fragment updates based on a single distributional structure, the Inverse G-Wishart distribution with $G \in \{G_{\text{full}}, G_{\text{diag}}\}$, are correct. The good news is that these algorithms only need to be derived once. Their implementations, within a suite of computer programmes for carrying out variational message passing for models containing variance and covariance matrix parameters, can be isolated into subroutines which, once working as intended, do not have to be revisited ever again. Given the quintessence of variance and covariance parameters in throughout statistics and machine learning, Algorithms 1 and 2 are important and fundamental contributions to variational message passing.

Acknowledgements

We are grateful to two referees for their comments and suggestions. This research was supported by Australian Research Council Discovery Project DP140100441.

References

- ARMAGAN, A., DUNSON, D.B. & CLYDE, M. (2011). Generalized beta mixtures of Gaussians. In *Advances in Neural Information Processing Systems 24*, eds. J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F. Pereira and K.Q. Weinberger, pp. 523–531. La Jolla, CA: The Neural Information Processing Systems (NIPS) Foundation.
- ASSAF, A.G., LI, G., SONG, H. & TSIONAS, M.G. (2019). Modeling and forecasting regional tourism demand using the Bayesian global vector autoregressive (BGVAR) model. *Journal of Travel Research* **58**, 383–397.
- ATAY-KAYIS, A. & MASSAM, H. (2005). A Monte Carlo method for computing the marginal likelihood in nondecomposable Gaussian graphical models. *Biometrika* **92**, 317–335.
- ATTIAS, H. (1999). Inferring parameters and structure of latent variable models by variational Bayes. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, eds. K.B. Laskey and H. Prade, pp. 21–30. San Francisco: Morgan Kaufmann.

- BISHOP, C.M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- CHEN, W.Y. & WAND, M.P. (2020). Factor graph fragmentation of expectation propagation. *Journal of the Korean Statistical Society* **49**, 722–756.
- CONTI, G., FRÜHWIRTH-SCHNATTER, S., HECKMAN, J.J. & PIATEK, R. (2014). Bayesian exploratory factor analysis. *Journal of Econometrics* **183**, 31–57.
- DAWID, A.P. & LAURITZEN, S.L. (1993). Hyper Markov laws in the statistical analysis of decomposable graphical models. *The Annals of Statistics* **21**, 1272–1317.
- GELMAN, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian Analysis* **1**, 515–534.
- GELMAN, A., CARLIN, J.B., STERN, H.S., DUNSON, D.B., VEHTARI, A. & RUBIN, D.B. (2014). *Bayesian Data Analysis*, 3rd edn. Boca Raton, Florida: CRC Press.
- GENTLE, J.E. (2007). *Matrix Algebra*. New York: Springer.
- HAREZLAK, J., RUPPERT, D. & WAND, M.P. (2018). *Semiparametric Regression with R*. New York: Springer.
- HUANG, A. & WAND, M.P. (2013). Simple marginally noninformative prior distributions for covariance matrices. *Bayesian Analysis* **8**, 439–452.
- LANGE, K.L., LITTLE, R.J. & TAYLOR, J.M. (1989). Robust statistical modeling using the t distribution. *Journal of the American Statistical Association* **84**, 881–896.
- LETAC, G. & MASSAM, H. (2007). Wishart distributions for decomposable graphs. *The Annals of Statistics* **35**, 1278–1323.
- MAESTRINI, L. & WAND, M.P. (2018). Variational message passing for skew t regression. *Stat* **7**, e196, 1–11.
- MAGNUS, J.R. & NEUDECKER, H. (2019). *Matrix Differential Calculus with Applications in Statistics and Econometrics*, 3rd edn. Chichester U.K.: John Wiley & Sons.
- MCCULLOCH, C.E., SEARLE, S.R. & NEUHAUS, J.M. (2008). *Generalized, Linear, and Mixed Models*, 2nd edn. New York: John Wiley & Sons.
- MCLEAN, M. & WAND, M. (2019). Variational message passing for elaborate response regression models. *Bayesian Analysis* **14**, 371–398.
- MINKA, T.P. (2001). Expectation propagation for approximate Bayesian inference. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, pp. 362–369. Burlington, Massachusetts: Morgan Kaufmann.
- MINKA, T.P. (2005). Divergence measures and message passing. *Microsoft Research Technical Report Series MSR-TR-2005-173*, 1–17.
- MUIRHEAD, R.J. (1982). *Aspects of Multivariate Statistical Theory*. New York: John Wiley & Sons.
- MULDER, J. & PERICCHI, L.R. (2018). The matrix- F prior for estimating and testing covariance matrices. *Bayesian Analysis* **13**, 1193–1214.
- NOLAN, T.H., MENICTAS, M. & WAND, M.P. (2020). Streamlined computing for variational inference with higher level random effects. *Journal of Machine Learning Research* **21**, 1–62.
- NOLAN, T.H. & WAND, M.P. (2017). Accurate logistic variational message passing: algebraic and numerical details. *Stat* **6**, 102–112.
- POLSON, N.G. & SCOTT, J.G. (2012). On the half-Cauchy prior for a global scale parameter. *Bayesian Analysis* **7**, 887–902.
- R CORE TEAM (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. Available from URL: <https://www.R-project.org/>.
- ROVERATO, A. (2000). Cholesky decomposition of a hyper inverse Wishart matrix. *Biometrika* **87**, 99–112.
- STAN DEVELOPMENT TEAM (2020). RStan: the R interface to Stan. Available from URL: <http://mc-stan.org/>. R package version 2.21.2.
- UHLER, C., LENKOSKI, A. & RICHARDS, D. (2018). Exact formulas for the normalizing constants of Wishart distributions for graphical models. *The Annals of Statistics* **46**, 90–118.
- WAND, M.P. (2017). Fast approximate inference for arbitrarily large semiparametric regression models via message passing. *Journal of the American Statistical Association* **112**, 137–168.
- WINN, J. & BISHOP, C.M. (2005). Variational message passing. *Journal of Machine Learning Research* **6**, 661–694.