

# The Inverse G-Wishart Distribution and Variational Message Passing

BY L. MAESTRINI AND M.P. WAND

*University of Technology Sydney*

15th June, 2020

## Abstract

Message passing on a factor graph is a powerful paradigm for the coding of approximate inference algorithms for arbitrarily graphical large models. The notion of a factor graph fragment allows for compartmentalization 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.

*Keywords:* 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 modularization 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 neighboring 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. Whilst 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) introduced the concept of *factor graph fragments*, or *fragments* for short, for compartmentalization of variational message passing into atom-like components. Chen & 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 Sections 2.4–2.5 of Wand (2017).

Sections 4.1.2–4.1.3 of Wand (2017) introduce 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 stan-

dard deviation parameters such as members of the Half- $t$  family (Gelman, 2006 ; Polson & Scott, 2012). The extension to the covariance matrix case, for which there is the option to impose Uniform distribution priors over the interval  $(-1, 1)$  on correlation parameters, is elucidated in Huang & Wand (2013). These two fragments arise in many classes of Bayesian models, such as both Gaussian and generalized response linear mixed models (e.g. McCulloch *et al.*, 2008), Bayesian factor models (e.g. Conti *et al.*, 2014), vector autoregressive models (e.g. Assaf *et al.*, 2019), and generalized additive mixed models and group-specific curve models (e.g. Harezlak *et al.*, 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, 2020) 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 family 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 Sections 4.1.2 and 4.1.3 of Wand (2017), 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 a 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.

Let  $G$  be an undirected graph with  $d$  nodes labeled  $1, \dots, d$  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  $G_{\text{full}}$ . 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  $G_{\text{diag}}$ .

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

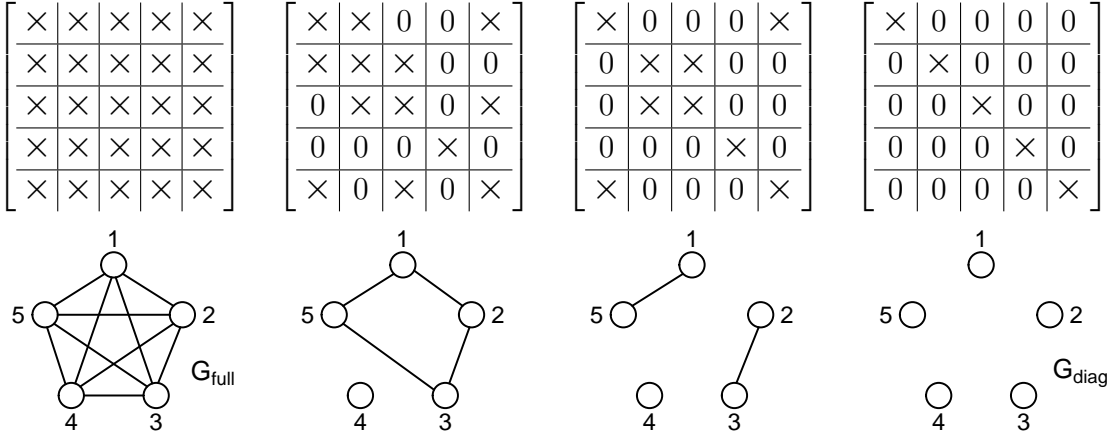


Figure 1: The zero/non-zero entries of four  $5 \times 5$  symmetric matrices with non-zero entries denoted by  $\times$ . Underneath each matrix is the 5-node undirected graph that the matrix respects. The nodes are numbered according to the rows and columns of the matrices. A graph edge is present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

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  $I(\mathcal{P})$  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  $x_j \stackrel{\text{ind.}}{\sim} \mathcal{D}_j$ ,  $1 \leq j \leq d$ . 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  $d^2 \times \{\frac{1}{2}d(d+1)\}$  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, 1999). 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  $\mathbf{D}_d^+ \equiv (\mathbf{D}_d^T \mathbf{D}_d)^{-1} \mathbf{D}_d^T$  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  $\delta$  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

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

Obtaining an expression for the normalizing factor of a general G-Wishart density function is a challenging problem and recently was resolved by Uhler *et al.* (2018). In the special case where  $G$  is a decomposable graph a relatively simple expression for the normalizing factor exists and is given, for example, by equation (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 normalizing 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 Section 3 of Atay-Kayis & Massam (2005), the G-Wishart distribution has connections with other distributional constructs such as the hyper Wishart law defined by Dawid & Lauritzen (1993).

Let  $G_{\text{full}}$  be the totally connected  $d$ -node undirected graph and  $G_{\text{diag}}$  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 normalizing 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*

$$\begin{aligned} \mathfrak{p}(\mathbf{X}) = & \frac{|\mathbf{\Lambda}|^{(\delta+d-1)/2}}{2^{d(\delta+d-1)/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma(\frac{\delta+d-j}{2})} |\mathbf{X}|^{(\delta-2)/2} \exp\{-\frac{1}{2}\text{tr}(\mathbf{\Lambda}\mathbf{X})\} \\ & \times I(\mathbf{X} \text{ a symmetric and positive definite } d \times d \text{ matrix}). \end{aligned} \quad (2)$$

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 the shape parameter used here is different from that commonly used for the Wishart distribution. For example, in Table A.1 of Gelman *et al.* (2014) the shape parameter is denoted by  $\nu$  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.

### 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 Chi-Squared distribution for a scalar random variable  $x$ .

**Definition 2.** *Let  $x$  be a random variable. For  $\delta > 0$  and  $\lambda > 0$  we say that  $x$  has a Chi-Squared distribution with shape parameter  $\delta$  and rate parameter  $\lambda$ , and write  $x \sim \chi^2(\delta, \lambda)$ , if and only if the density function of  $x$  satisfies*

$$\mathfrak{p}(x) = \frac{(\lambda/2)^{\delta/2}}{\Gamma(\delta/2)} x^{(\delta-2)/2} \exp(-\frac{1}{2}\lambda x) I(x > 0).$$

The G-Wishart( $G_{\text{diag}}, \delta, \mathbf{\Lambda}$ ) distribution is tied intimately to the Chi-Squared 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

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

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

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

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.

1. 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  $\mathbf{X} \sim \text{G-Wishart}(G_{\text{diag}}, \delta, \text{diag}\{\text{diagonal}(\mathbf{\Lambda})\})$ .
2. 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 Chi-Squared random variables with shape parameter  $\delta$  and rate parameters equalling the diagonal entries of  $\mathbf{\Lambda}$ .
3. 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

$$p(\mathbf{X}) \propto \exp \left\{ \begin{bmatrix} \log |\mathbf{X}| \\ \text{vech}(\mathbf{X}) \end{bmatrix}^T \begin{bmatrix} \frac{1}{2}(\delta - 2) \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\mathbf{\Lambda}) \end{bmatrix} \right\} = \exp\{\mathbf{T}(\mathbf{X})^T \boldsymbol{\eta}\} \quad (3)$$

where

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

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

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

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  $p_{\mathbf{X}}$  and  $p_{\mathbf{Y}}$  respectively. Then the density function of  $\mathbf{Y}$  is

$$p_{\mathbf{Y}}(\mathbf{Y}) = p_{\mathbf{X}}(\mathbf{Y}^{-1}) |J(\mathbf{Y})| \quad (4)$$

where

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

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 (4). 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

$$\text{if } G = G_{\text{diag}} \text{ then } |J(\mathbf{Y})| = |\mathbf{Y}|^{-2} \text{ for any } d \in \mathbb{N}. \quad (5)$$

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

$$\text{if } G = G_{\text{full}} \text{ then } |J(\mathbf{Y})| = |\mathbf{Y}|^{-(d+1)}. \quad (6)$$

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

Combining (4), (5) and (6) we have:

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

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

Whilst 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 3.** 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  $\mathbf{\Lambda}$ .

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

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

(b) If  $G = G_{\text{diag}}$  then say that  $\mathbf{X}$  has an Inverse G-Wishart distribution with graph  $G$ , shape parameter  $\xi$  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

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

(c) 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  $\xi$  used in Definition 3 is a reasonable compromise between various competing parameterization 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  $\xi$  matches the shape parameter in the most common parameterization of the Inverse Chi-Squared distribution such as that used in Table A.1 of Gelman *et al.* (2014).

In 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

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

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, Table A.1 of Gelman *et al.* (2014).

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

**Definition 4.** 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  $\delta$  and rate parameter  $\lambda$ , 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

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

We are now ready to state:

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

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

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  $p(\mathbf{X}) > 0$ ,

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

where

$$\mathbf{T}(\mathbf{X}) \equiv \left[\begin{array}{c} \log|\mathbf{X}| \\ \text{vech}(\mathbf{X}^{-1}) \end{array}\right] \quad \text{and} \quad \boldsymbol{\eta} \equiv \left[\begin{array}{c} \eta_1 \\ \boldsymbol{\eta}_2 \end{array}\right] = \left[\begin{array}{c} -\frac{1}{2}(\xi+2) \\ -\frac{1}{2}\mathbf{D}_d^T \text{vec}(\mathbf{\Lambda}) \end{array}\right] \quad (7)$$

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

$$\begin{cases} \xi &= -2\eta_1 - 2, \\ \mathbf{\Lambda} &= -2 \text{vec}^{-1}(\mathbf{D}_d^{+T} \boldsymbol{\eta}_2). \end{cases} \quad (8)$$

As explained in Section S.1 of the web-supplement, alternatives to (7) 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

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

### 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 Connection with 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( $G, \xi, \mathbf{\Lambda}$ ) 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  $\Sigma$  is a  $d \times d$  symmetric positive definite matrix then, for  $\nu > 0$  and  $s_1, \dots, s_d > 0$ , the specification

$$\begin{aligned} \Sigma | \mathbf{A} &\sim \text{Inverse-G-Wishart}\left(G_{\text{full}}, \nu + 2d - 2, \mathbf{A}^{-1}\right), \\ \mathbf{A} &\sim \text{Inverse-G-Wishart}\left(G_{\text{diag}}, 1, \{\nu \text{diag}(s_1^2, \dots, s_d^2)\}^{-1}\right) \end{aligned} \tag{9}$$

places a Huang-Wand distribution on  $\Sigma$  with shape parameter  $\nu$  and scale parameters  $s_1, \dots, s_d$ .

The specification (9) matches (2) of Huang & Wand (2013) but with some differences in notation. Firstly,  $d$  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  $\mathbf{A}$  via the expression  $\text{diag}(a_1, \dots, a_d) = 2\nu \mathbf{A}$ .

As discussed in Huang & Wand (2013), special cases of (9) correspond to marginally noninformative prior specification of the covariance matrix  $\Sigma$  in the sense that the standard deviation parameters  $\sigma_j \equiv (\Sigma)_{jj}^{1/2}$ ,  $1 \leq j \leq d$ , 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  $\rho_{jj'} \equiv (\Sigma)_{jj'}^{1/2} / (\sigma_j \sigma_{j'})$ , have a Uniform distribution over the interval  $(-1, 1)$  when  $\nu = 2$ . 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

$$\Sigma \sim \text{Huang-Wand}(s_1, \dots, s_d) \tag{10}$$

as a shorthand for (9) with  $\nu = 2$ .

### 4 Variational Message Passing Background

The overarching goal of this article is 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 variational message passing here. Instead we refer the reader to Sections 2–4 of Wand (2017) 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  $\theta$  denotes a generic stochastic variable that is a neighbor of  $f$  in the factor graph then the message passed from  $f$  to  $\theta$  and the message passed from  $\theta$  to  $f$  are both functions of  $\theta$  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

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

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

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

## 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  $\Theta$ :

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

for a  $d$ -node undirected graph  $G_{\Theta}$ , scalar shape parameter  $\xi_{\Theta}$  and scale matrix  $\Lambda_{\Theta}$ . The fragment's factor is

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



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

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

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

Therefore, the natural parameter update is

$$\boldsymbol{\eta}_{p(\Theta) \rightarrow \Theta} \leftarrow \begin{bmatrix} -\frac{1}{2}(\xi_{\Theta} + 2) \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\Lambda_{\Theta}) \end{bmatrix}.$$

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:**

$$\boldsymbol{\eta}_{p(\Theta) \rightarrow \Theta} \leftarrow \begin{bmatrix} -\frac{1}{2}(\xi_{\Theta} + 2) \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\Lambda_{\Theta}) \end{bmatrix}; \quad G_{p(\Theta) \rightarrow \Theta} \leftarrow G_{\Theta}$$

**Outputs:**  $G_{p(\Theta) \rightarrow \Theta}, \boldsymbol{\eta}_{p(\Theta) \rightarrow \Theta}$ .

---

## 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  $\boldsymbol{\Sigma}$ :

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

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

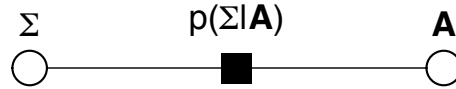


Figure 3: *Diagram of the iterated Inverse G-Wishart fragment.*

The factor of the iterated Inverse G-Wishart fragment is, as a function of both  $\boldsymbol{\Sigma}$  and  $\mathbf{A}$ ,

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

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

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

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  $\boldsymbol{\Sigma}$  and  $\mathbf{A}$  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  $\boldsymbol{\Sigma}$  and  $p(\boldsymbol{\Sigma} | \mathbf{A})$  have the same graph. The Inverse G-Wishart messages passed between  $\mathbf{A}$  and  $p(\boldsymbol{\Sigma} | \mathbf{A})$  have the same graph.

Under these constraints, and in view of e.g. (7) of Wand (2017), the message passed from  $\boldsymbol{\Sigma}$  to  $p(\boldsymbol{\Sigma} | \mathbf{A})$  has the form

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

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

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

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:**  $G_{\mathbf{A} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})} \in \{G_{\text{full}}, G_{\text{diag}}\}$ .

**Natural Parameter Inputs:**  $\boldsymbol{\eta}_{\boldsymbol{\Sigma} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}$ ,  $\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}$ ,  $\boldsymbol{\eta}_{\mathbf{A} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}$ ,  $\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}}$ .

**Updates:**

$$G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} \leftarrow G \ ; \ G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} \leftarrow G_{\mathbf{A} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}$$

$$\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \boldsymbol{\Sigma}} \leftarrow \boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} + \boldsymbol{\eta}_{\boldsymbol{\Sigma} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}$$

$$\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}} \leftarrow \boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} + \boldsymbol{\eta}_{\mathbf{A} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{full}} \text{ then } \omega_1 \leftarrow (d+1)/2$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{diag}} \text{ then } \omega_1 \leftarrow 1$$

$$E_{\mathbf{q}}(\mathbf{A}^{-1}) \leftarrow \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_1 + \omega_1 \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_2 \right) \right\}^{-1}$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} = G_{\text{diag}} \text{ then } E_{\mathbf{q}}(\mathbf{A}^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left( E_{\mathbf{q}}(\mathbf{A}^{-1}) \right) \right\}$$

$$\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} \leftarrow \begin{bmatrix} -\frac{1}{2} (\xi + 2) \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec} \left( E_{\mathbf{q}}(\mathbf{A}^{-1}) \right) \end{bmatrix}$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} = G_{\text{full}} \text{ then } \omega_2 \leftarrow (d+1)/2$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} = G_{\text{diag}} \text{ then } \omega_2 \leftarrow 1$$

$$E_{\mathbf{q}}(\boldsymbol{\Sigma}^{-1}) \leftarrow \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \boldsymbol{\Sigma}})_1 + \omega_2 \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \boldsymbol{\Sigma}})_2 \right) \right\}^{-1}$$

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{diag}} \text{ then } E_{\mathbf{q}}(\boldsymbol{\Sigma}^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left( E_{\mathbf{q}}(\boldsymbol{\Sigma}^{-1}) \right) \right\}$$

$$\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} \leftarrow \begin{bmatrix} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec} \left( E_{\mathbf{q}}(\boldsymbol{\Sigma}^{-1}) \right) \end{bmatrix}$$

**Outputs:**  $G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}$ ,  $G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}}$ ,  $\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}$ ,  $\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}}$ .

---

## 6.1 Corrections to Section 4.1.3 of Wand (2017)

The iterated Inverse G-Wishart fragment was introduced in Section 4.1.3 of Wand (2017) 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 Section 4.1.3 of Wand (2017).

Firstly, in Wand (2017)  $\Theta_1$  plays the role of  $\Sigma$  and  $\Theta_2$  plays the role of  $A$ . The dimension of  $\Theta_1$  and  $\Theta_2$  is denoted by  $d^\Theta$ . The first displayed equation of Section 4.1.3 is

$$\Theta_1 | \Theta_2 \sim \text{Inverse-G-Wishart}(G, \kappa, \Theta_2^{-1}) \quad (11)$$

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, (11) involves a different parameterization to that used throughout this article. Therefore, our first correction is to replace the first displayed equation of Section 4.1.3 of Wand (2017) by:

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

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

The following sentence in Section 4.1.3 of Wand (2017): “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 equation (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 parameterization of the Inverse G-Wishart distribution, as used in the current article, equation (32) of Wand (2017) should be

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

The equation in Section 4.1.3 of Wand (2017):

$$“\boldsymbol{\eta}_{p(\Theta_1 | \Theta_2) \rightarrow \Theta_2} \leftarrow \begin{bmatrix} -\kappa/2 \\ -\frac{1}{2} \text{vec} \left( E_{p(\Theta_1 | \Theta_2) \rightarrow \Theta_2}(\Theta_1^{-1}) \right) \end{bmatrix} ”$$

should be replaced by

$$“\boldsymbol{\eta}_{p(\Theta_1 | \Theta_2) \rightarrow \Theta_2} \leftarrow \begin{bmatrix} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2} \text{vec} \left( E_{p(\Theta_1 | \Theta_2) \rightarrow \Theta_2}(\Theta_1^{-1}) \right) \end{bmatrix} ”$$

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

Lastly the iterated Inverse G-Wishart fragment natural parameter updates given by equations (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  $\sigma^2$  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  $\sigma^2$  be a variance parameter and consider the prior imposition

$$\sigma^2 \sim \text{Inverse-Gamma}(\alpha_{\sigma^2}, \beta_{\sigma^2}) \tag{12}$$

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

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

Note that the Inverse Chi-Squared and Inverse Gamma distributions are simple reparameterizations 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 (12) 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  $\mathbf{X}$  is defined to have an Inverse Wishart distribution with shape parameter  $\kappa$  and scale matrix  $\Sigma$ , written  $\mathbf{X} \sim \text{Inverse-Wishart}(\kappa, \Lambda)$ , if and only if the density function of  $\mathbf{X}$  is

$$\begin{aligned} p(\mathbf{X}) = & \frac{|\Lambda|^{\kappa/2}}{2^{\kappa d/2} \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma\left(\frac{\kappa+1-j}{2}\right)} |\mathbf{X}|^{-(\kappa+d+1)/2} \exp\left\{-\frac{1}{2}\text{tr}(\Lambda\mathbf{X}^{-1})\right\} \\ & \times I(\mathbf{X} \text{ a symmetric and positive definite } d \times d \text{ matrix}). \end{aligned} \tag{13}$$

Note that this is the common parameterization of the Inverse Wishart distribution (e.g. Table A.1 of Gelman *et al.*, 2014). Crucially, (13) uses a *different* shape parametrization from that used for the Inverse G-Wishart distribution in Definition 3 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 parameterization as given in (13), is more common when imposing a prior on a covariance matrix.

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

$$\Sigma \sim \text{Inverse-Wishart}(\kappa_\Sigma, \Lambda_\Sigma) \quad (14)$$

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

$$\sigma \sim \text{Half-}t(s_\sigma, \nu_\sigma) \quad (15)$$

for a scale parameter  $s_\sigma > 0$  and a degrees of freedom parameter  $\nu_\sigma > 0$ . The density function corresponding to (15) is such that  $p(\sigma) \propto \{1 + (\sigma/s_\sigma)^2/\nu_\sigma\}^{-(\nu_\sigma+1)/2} I(\sigma > 0)$ . This is equivalent to

$$\sigma^2|a \sim \text{Inverse-}\chi^2(\nu_\sigma, 1/a) \quad \text{and} \quad a \sim \text{Inverse-}\chi^2(1, 1/s_\sigma^2). \quad (16)$$

Since  $d = 1$ , the graphs  $G_{\text{full}}$  and  $G_{\text{diag}}$  are the same – a single node graph. Treating  $\sigma^2$  and  $a$  as  $1 \times 1$  matrices we can re-write (16) as

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

The specification

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

involves calling Algorithm 1 with

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

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

$$\boldsymbol{\eta}_{p(\Theta) \rightarrow \Theta} = \boldsymbol{\eta}_{p(a) \rightarrow a}.$$

The specification

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

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

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

and

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

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

##### 7.4.1 The Half-Cauchy Special Case

The special case of

$$\sigma \sim \text{Half-Cauchy}(s_\sigma). \quad (17)$$

corresponds to  $\nu_\sigma = 1$ . The density function corresponding to (17) is such that  $p(\sigma) \propto \{1 + (\sigma/s_\sigma)^2\}^{-1} I(\sigma > 0)$ . 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 (10)

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

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 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.

| prior specification  | Algorithm 1       |                           |   | Algorithm 2    |                   |   |
|--|-------------------|---------------------------|---|----------------|-------------------|---|
|  | $G_{\Theta}$      | $\xi_{\Theta}$            | $\Lambda_{\Theta}$  | $\xi$          | $G$               | $G_{A \rightarrow \mathbf{p}(\Sigma \mathbf{A})}$ |
| $\sigma^2 \sim \text{Inverse-}\chi^2(\delta_{\sigma^2}, \lambda_{\sigma^2})$ | $G_{\text{full}}$ | $\delta_{\sigma^2}$       | $\lambda_{\sigma^2}$  | N.A.           | N.A.              | N.A.  |
| $\sigma^2 \sim \text{Inv.-Gamma}(\alpha_{\sigma^2}, \beta_{\sigma^2})$       | $G_{\text{full}}$ | $2\alpha_{\sigma^2}$      | $2\beta_{\sigma^2}$   | N.A.           | N.A.              | N.A.  |
| $\Sigma \sim \text{Inv.-Wishart}(\kappa_{\Sigma}, \Lambda_{\Sigma})$         | $G_{\text{full}}$ | $\kappa_{\Sigma} + d - 1$ | $\Lambda_{\Sigma}$  | N.A.           | N.A.              | N.A.  |
| $\sigma \sim \text{Half-}t(s_{\sigma}, \nu_{\sigma})$                        | $G_{\text{diag}}$ | 1                         | $(\nu_{\sigma} s_{\sigma}^2)^{-1}$                              | $\nu_{\sigma}$ | $G_{\text{full}}$ | $G_{\text{diag}}$                                 |
| $\sigma \sim \text{Half-Cauchy}(s_{\sigma})$                                 | $G_{\text{diag}}$ | 1                         | $(s_{\sigma}^2)^{-1}$   | 1              | $G_{\text{full}}$ | $G_{\text{diag}}$                                 |
| $\Sigma \sim \text{Huang-Wand}$<br>$(s_{\Sigma,1}, \dots, s_{\Sigma,d})$     | $G_{\text{diag}}$ | 1                         | $\{2 \text{diag}(s_{\Sigma,1}^2, \dots, s_{\Sigma,d}^2)\}^{-1}$ | $2d$           | $G_{\text{full}}$ | $G_{\text{diag}}$                                 |

Table 1: Specifications of inputs of Algorithms 1 and 2 for several variance, standard deviation and covariance matrix prior impositions. The abbreviation N.A. stands for not applicable since Algorithm 2 is not needed for the first three prior impositions.

## 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 *et al.*, 1989). The notation  $y \sim t(\mu, \sigma, \nu)$  indicates that the random variable  $y$  has a  $t$  distribution with location parameter  $\mu$ , scale parameter  $\sigma > 0$  and degrees of freedom parameter  $\nu > 0$ . The corresponding density function of  $y$  is

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

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

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



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

$$\mathbf{y} \equiv \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{bmatrix} \quad \text{and} \quad \mathbf{Z} \equiv \text{blockdiag}(\mathbf{Z}_i)_{1 \leq i \leq m} \quad (18)$$

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  $y_\ell$  be the  $\ell$ th entry of  $\mathbf{y}$ ,  $1 \leq \ell \leq N$ . The family of Bayesian  $t$  response linear mixed models that we consider is

$$\begin{aligned} y_\ell | \boldsymbol{\beta}, \mathbf{u}, \sigma &\stackrel{\text{ind.}}{\sim} t((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_\ell, \sigma, \nu), \quad 1 \leq \ell \leq N, \quad \mathbf{u} | \boldsymbol{\Sigma} \sim N(\mathbf{0}, \mathbf{I}_m \otimes \boldsymbol{\Sigma}), \\ \boldsymbol{\beta} &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}), \quad \sigma \sim \text{Half-Cauchy}(s_\sigma), \quad \frac{1}{2}\nu \sim \text{Moon-Rock}(0, \lambda_\nu), \\ \boldsymbol{\Sigma} &\sim \text{Huang-Wand}(s_{\boldsymbol{\Sigma},1}, \dots, s_{\boldsymbol{\Sigma},q}) \end{aligned} \quad (19)$$

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  $p(x) \propto \{x^x/\Gamma(x)\}^\alpha \exp(-\beta x)I(x > 0)$ . In the variational message passing treatment of the degrees of freedom parameter it is simpler to work with

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

After the approximate posterior density function of  $v$  is obtained via variational message passing, it is trivial to then obtain the same for  $\nu$ . Hence, we work with  $v$ , rather than  $\nu$ , in the upcoming description of variational message passing-based fitting and inference for (19).

Next note that

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

is equivalent to

$$\mathbf{y}_\ell | \boldsymbol{\beta}, \mathbf{u}, \sigma^2, b_\ell \sim N((\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u})_\ell, b_\ell \sigma^2), \quad b_\ell | v \stackrel{\text{ind.}}{\sim} \text{Inverse-}\chi^2(2v, 2v), \quad (20)$$

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

$$\sigma^2 | a \sim \text{Inverse-}\chi^2(1, 1/a) \quad \text{and} \quad a \sim \text{Inverse-}\chi^2(1, 1/s_\sigma^2) \quad (21)$$

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

$$\begin{aligned} \boldsymbol{\Sigma} | \mathbf{A} &\sim \text{Inverse-G-Wishart}(G_{\text{full}}, 2q, \mathbf{A}^{-1}), \\ \mathbf{A} &\sim \text{Inverse-G-Wishart}\left(G_{\text{diag}}, 1, \{2 \text{diag}(s_{\boldsymbol{\Sigma},1}^2, \dots, s_{\boldsymbol{\Sigma},q}^2)\}^{-1}\right). \end{aligned} \quad (22)$$

Substitution of (20), (21) and (22) into (19) leads to the hierarchical Bayesian model depicted as a directed acyclic graph in Figure 4 with  $\mathbf{b} \equiv (b_1, \dots, b_N)$ . 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

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b} | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u}, a, \mathbf{A}, \mathbf{b}) q(\sigma^2, \boldsymbol{\Sigma}, v) \quad (23)$$

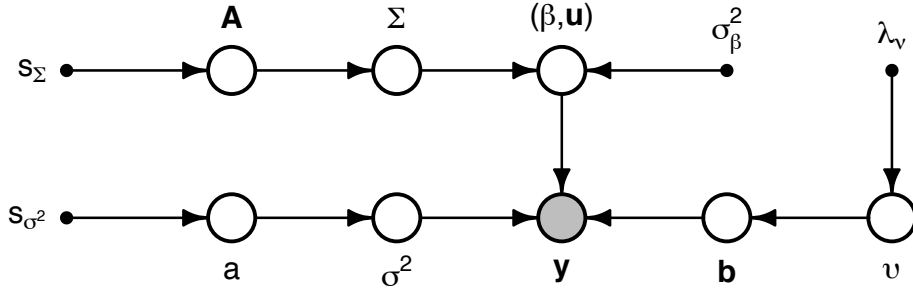


Figure 4: Directed acyclic graph corresponding to the  $t$  response linear mixed model (19) with auxiliary variable representations (20)–(22). The shaded circle corresponds to the observed data. The unshaded circles correspond to model parameters and auxiliary variables. The small solid circles correspond to hyperparameters.

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 factorizations

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

and so the restriction given in (23) is equivalent to

$$p(\boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b} | \mathbf{y}) \approx q(\boldsymbol{\beta}, \mathbf{u})q(\sigma^2)q(v)q(\boldsymbol{\Sigma})q(a)q(\mathbf{A}) \prod_{\ell=1}^N q(b_\ell). \quad (24)$$

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 factorization given in (24) and filled-in rectangles corresponding to factors on the right-hand side of

$$p(\mathbf{y}, \boldsymbol{\beta}, \mathbf{u}, \sigma^2, v, \boldsymbol{\Sigma}, a, \mathbf{A}, \mathbf{b}) = p(\mathbf{A})p(a)p(\boldsymbol{\Sigma} | \mathbf{A})p(\sigma^2 | a)p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma})p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{u}, \sigma^2, \mathbf{b})p(\mathbf{b} | v)p(v). \quad (25)$$

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 neighboring 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 neighboring stochastic nodes to the factor.
3. Update the parameter vectors of the messages passed from the factor to its neighboring 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:

$$\begin{aligned} & \text{the message passed from } \boldsymbol{\Sigma} \text{ to } p(\boldsymbol{\beta}, \mathbf{u} | \boldsymbol{\Sigma}) \\ & = \text{the message passed from } p(\boldsymbol{\Sigma} | \mathbf{A}) \text{ to } \boldsymbol{\Sigma} \text{ in the previous iteration.} \end{aligned} \quad (26)$$

In terms of natural parameter vector updates, (26) corresponds to:

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

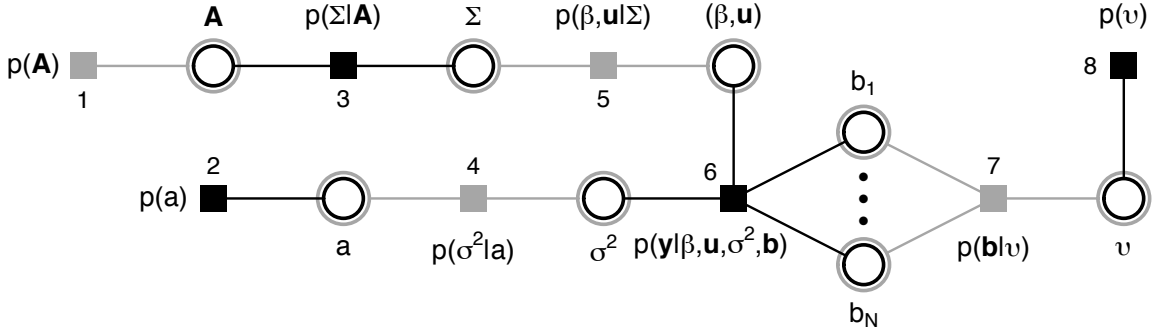


Figure 5: Factor graph corresponding to the  $t$  response linear mixed model (19) with auxiliary variable representations (20)–(22). The circular nodes correspond to stochastic nodes in the  $q$ -density factorization in (24). The rectangular nodes correspond to the factors on the right-hand side of (25). The fragments are numbered 1 to 8 according to appearance from left to right. Shading is used to show the distinction between adjacent fragments.

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, equation (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

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

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 penalization 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 Algorithm 2 of McLean & Wand (2019).
- 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  $\theta$  is

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

This corresponds to the prior density function having exponential family form

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

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

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

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

To demonstrate variational message passing for fitting and inference for model (19), we simulated data according to the dimension values  $p = q = 2$  and the true parameter values

$$\boldsymbol{\beta}_{\text{true}} = \left[ \begin{array}{c} -0.58 \\ 1.89 \end{array} \right], \quad \sigma_{\text{true}}^2 = 0.2, \quad \boldsymbol{\Sigma}_{\text{true}} = \left[ \begin{array}{cc} 2.58 & 0.22 \\ 0.22 & 1.73 \end{array} \right] \quad \text{and} \quad \nu_{\text{true}} = 1.5. \quad (27)$$

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

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

We ran the variational message passing algorithm as described above until the relative change the variational parameters was below  $10^{-10}$ . as well as Markov chain Monte Carlo via the R language (R Core Team, 2020) package `rstan` (Stan Development Team, 2019). For Markov chain Monte Carlo fitting, a warmup of size 1000 was used, followed by chains of size 5000, thinned by a factor of 5, 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  $\boldsymbol{\Sigma}$  matrix, according to the notation  $(\boldsymbol{\Sigma})_{11} = \sigma_1$ ,  $(\boldsymbol{\Sigma})_{22} = \sigma_2$  and  $(\boldsymbol{\Sigma})_{12} = \sigma_1 \sigma_2 \rho$ . For most of the stochastic nodes, the accuracy of variational message passing is seen to be very good. For  $\sigma$  and  $\nu$ , some under-approximation of the spread is apparent.

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,

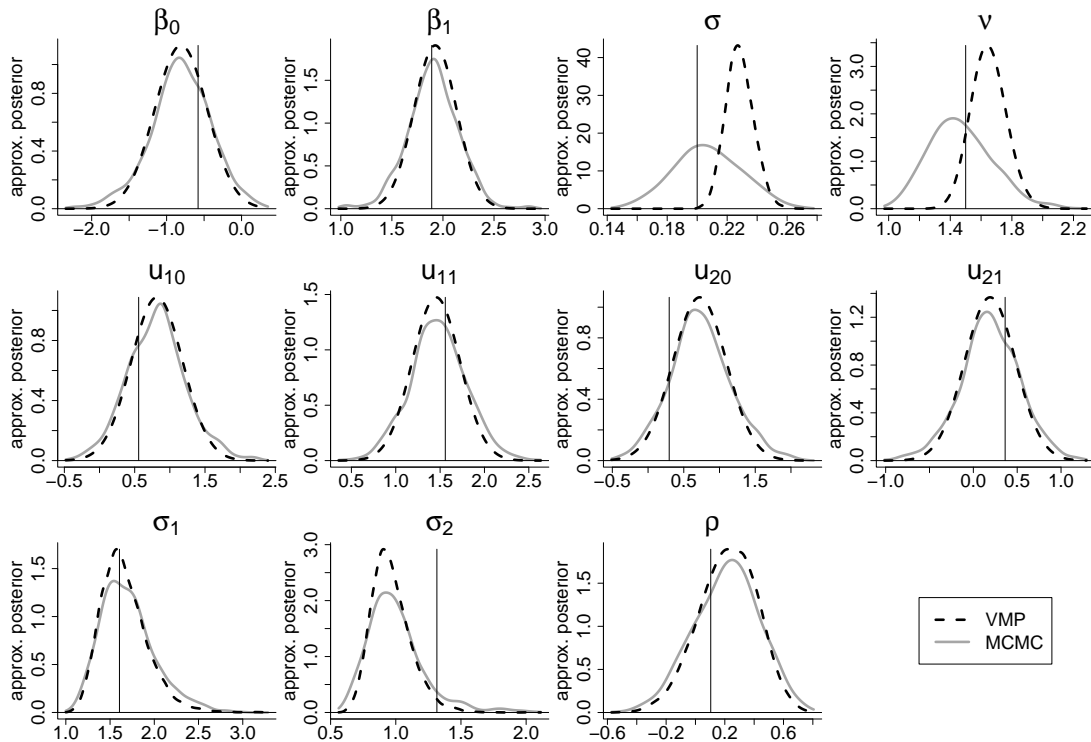


Figure 6: Approximate posterior density functions for the parameters in model (19) based on both variational message passing (VMP) and Markov chain Monte Carlo (MCMC) algorithms applied to data simulated according to the true values (27) and sample sizes, predictor values and hyperparameter values as described in the text. The vertical lines indicate true parameter values.

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 Algorithm 2 are important and fundamental contributions to variational message passing.

## Acknowledgements

This research was supported by Australian Research Council Discovery Project DP140100441.

## References

- 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.
- Attias, H. (1999). Inferring parameters and structure of latent variable models by variational Bayes. In Laskey, K.B. and Prade, H. (editors) *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, pp. 21–30. San Francisco: Morgan Kauffmann.

- Atay-Kayis, A. & Massam, H. (2005). A Monte Carlo method for computing marginal likelihood in nondecomposable Gaussian graphical models. *Biometrika*, **92**, 317–335.
- Bishop, C.M. (2006). *Pattern Recognition and Machine Learning*. New York: Springer.
- Chen, W.Y and Wand, M.P. (2020). Factor graph fragmentation of expectation propagation. *Journal of the Korean Statistical Society*, in press.
- 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. *Bayesian Analysis*, **1**, 515–533.
- Gelman, A., Carlin, J.B., Stern, H.S., Dunson, D.B., Vehtari, A. & Rubin, D.B. (2014). *Bayesian Data Analysis, Third Edition*, 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.A. and Taylor, J.M.G. (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.
- Magnus, J.R. & Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics, Revised Edition*. Chichester U.K.: Wiley
- McCulloch, C.E., Searle, S.R. & Neuhaus, J.M. (2008). *Generalized, Linear, and Mixed Models, Second Edition*. New York: John Wiley & Sons.
- McLean, M.W. & Wand, M.P. (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 J.S. Breese & D. Koller (eds), *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, pp. 362–369. Burlington, Massachusetts: Morgan Kaufmann.
- Minka, T. (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.

- Nolan, T.H., Menictas, M. and Wand, M.P. (2020). Streamlined computing for variational inference with higher level random effects. Unpublished manuscript available at <https://arxiv.org/abs/1903.06616>.
- Nolan, T.H. and 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 (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>
- Roverato, A. (2000). Cholesky decomposition of a hyper inverse Wishart matrix. *Biometrika*, **87**, 99–112.
- Stan Development Team (2019). RStan: the R interface to Stan. R package version 2.19.2. <http://mc-stan.org/>.
- Uhler, C., Lenkoski, A. and 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 (with discussion). *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.

Web-supplement for:

# The Inverse G-Wishart Distribution and Variational Message Passing

BY L. MAESTRINI AND M.P. WAND

*University of Technology Sydney*

## S.1 Natural Parameter Versions and Mappings

Throughout this article we use the “vech” versions of the natural parameter forms of the Multivariate Normal and Inverse G-Wishart distributions. However, Wand (2017) and McLean & Wand (2019) used “vec” versions of these distributions. The “vech” version has the attraction of being more compact since entries of symmetric matrices are not duplicated. However, adoption of the “vech” version entails use of duplication matrices. For implementation in the R language (R Core Team, 2020) we note that the function `duplication.matrix()` in the package `matrixcalc` (Novomestky, 2012) returns the duplication matrix of a given order.

First we explain the two versions for the Multivariate Normal distribution. Suppose that the  $d \times 1$  random vector  $\mathbf{v}$  has a  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution. Then the density function of  $\mathbf{v}$  is

$$p(\mathbf{v}) \propto \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vec}(\mathbf{v}\mathbf{v}^T) \end{bmatrix}^T \boldsymbol{\eta}_v^{\text{vec}} \right\} = \exp \left\{ \begin{bmatrix} \mathbf{v} \\ \text{vech}(\mathbf{v}\mathbf{v}^T) \end{bmatrix}^T \boldsymbol{\eta}_v^{\text{vech}} \right\}$$

where

$$\boldsymbol{\eta}_v^{\text{vec}} \equiv \begin{bmatrix} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ -\frac{1}{2} \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta}_v^{\text{vech}} \equiv \begin{bmatrix} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}.$$

The two natural parameter vectors can be mapped between each other using

$$\boldsymbol{\eta}_v^{\text{vech}} = \text{blockdiag}(\mathbf{I}_d, \mathbf{D}_d^T) \boldsymbol{\eta}_v^{\text{vec}} \quad \text{and} \quad \boldsymbol{\eta}_v^{\text{vec}} = \text{blockdiag}(\mathbf{I}_d, \mathbf{D}_d^{+T}) \boldsymbol{\eta}_v^{\text{vech}}. \quad (\text{S.1})$$

Now we explain the interplay between the “vec” and “vech” forms of the Inverse G-Wishart distribution. Let the  $d \times d$  matrix  $\mathbf{V}$  have an Inverse-G-Wishart( $G, \xi, \boldsymbol{\Lambda}$ ) distribution. Then the density function of  $\mathbf{V}$  is

$$p(\mathbf{V}) \propto \exp \left\{ \begin{bmatrix} \log |\mathbf{V}| \\ \text{vec}(\mathbf{V}^{-1}) \end{bmatrix}^T \boldsymbol{\eta}_V^{\text{vec}} \right\} = \exp \left\{ \begin{bmatrix} \log |\mathbf{V}| \\ \text{vech}(\mathbf{V}^{-1}) \end{bmatrix}^T \boldsymbol{\eta}_V^{\text{vech}} \right\}$$

where

$$\boldsymbol{\eta}_V^{\text{vec}} \equiv \begin{bmatrix} -\frac{1}{2}(\xi + 1) \\ -\frac{1}{2} \text{vec}(\boldsymbol{\Lambda}) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta}_V^{\text{vech}} \equiv \begin{bmatrix} -\frac{1}{2}(\xi + 1) \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec}(\boldsymbol{\Lambda}) \end{bmatrix}.$$

Mappings between the two natural parameter vectors are as follows:

$$\boldsymbol{\eta}_V^{\text{vech}} = \text{blockdiag}(1, \mathbf{D}_d^T) \boldsymbol{\eta}_V^{\text{vec}} \quad \text{and} \quad \boldsymbol{\eta}_V^{\text{vec}} = \text{blockdiag}(1, \mathbf{D}_d^{+T}) \boldsymbol{\eta}_V^{\text{vech}}. \quad (\text{S.2})$$

## S.2 Justification of Algorithm 2

We now provide justification for Algorithm 2, which is concerned with the graph and natural parameter updates for the iterated Inverse G-Wishart fragment.



### S.2.1 The Updates for $m_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$

As a function of  $\boldsymbol{\Sigma}$ ,

$$\log \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) = \begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vech}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^T \begin{bmatrix} -\frac{1}{2}(\xi + 2) \\ -\frac{1}{2}\mathbf{D}_d^T \text{vec}(\mathbf{A}^{-1}) \end{bmatrix} + \text{const}$$

where ‘const’ denotes terms that do not depend on  $\boldsymbol{\Sigma}$ . Hence

$$m_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) = \exp \left\{ \begin{bmatrix} \log |\boldsymbol{\Sigma}| \\ \text{vech}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix}^T \boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} \right\}$$

where

$$\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} = \begin{bmatrix} -\frac{1}{2}(\xi + 2) \\ -\frac{1}{2}\mathbf{D}_d^T \text{vec}(E_{\mathbf{q}}(\mathbf{A}^{-1})) \end{bmatrix} \quad (\text{S.3})$$

and  $E_{\mathbf{q}}$  denotes expectation with respect to the normalization of

$$m_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}}(\mathbf{A}) m_{\mathbf{A} \rightarrow \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A})}(\mathbf{A}).$$

Let  $q(\mathbf{A})$  denote this normalized density function. Then  $q(\mathbf{A})$  is an Inverse-G-Wishart distribution with graph  $G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} \in \{G_{\text{full}}, G_{\text{diag}}\}$  and natural parameter vector  $\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}}$ . From Result 6,

$$E_{\mathbf{q}}(\mathbf{A}^{-1}) = \begin{cases} \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_1 + \frac{1}{2}(d+1) \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_2 \right) \right\}^{-1} & \text{if } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{full}}, \\ \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_1 + 1 \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_2 \right) \right\}^{-1} & \text{if } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{diag}}. \end{cases}$$

Noting that the first factor of  $E_{\mathbf{q}}(\mathbf{A}^{-1})$  is  $(\boldsymbol{\eta}_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \leftrightarrow \mathbf{A}})_1 + \omega_1$ , where

$$\omega_1 = \omega_1(d, G) = \begin{cases} (d+1)/2 & \text{if } G = G_{\text{full}} \\ 1 & \text{if } G = G_{\text{diag}}, \end{cases}$$

the first update of  $E_{\mathbf{q}}(\mathbf{A}^{-1})$  in Algorithm 2 is justified. Lastly, we need to possibly adjust for the fact that  $m_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}}(\boldsymbol{\Sigma})$  is proportional to an Inverse G-Wishart density function with  $G = G_{\text{diag}}$ . This is achieved by the conditional step:

$$\text{If } G_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \boldsymbol{\Sigma}} = G_{\text{diag}} \text{ then } E_{\mathbf{q}}(\mathbf{A}^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left( E_{\mathbf{q}}(\mathbf{A}^{-1}) \right) \right\}.$$

### S.2.2 The Updates for $m_{\mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) \rightarrow \mathbf{A}}(\mathbf{A})$

As a function of  $\mathbf{A}$ ,

$$\log \mathbf{p}(\boldsymbol{\Sigma}|\mathbf{A}) = \begin{bmatrix} \log |\mathbf{A}| \\ \text{vech}(\mathbf{A}^{-1}) \end{bmatrix}^T \begin{bmatrix} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2}\mathbf{D}_d^T \text{vec}(\boldsymbol{\Sigma}^{-1}) \end{bmatrix} + \text{const}$$

where

$$\omega_2 = \omega_2(d, G) = \begin{cases} (d+1)/2 & \text{if } G = G_{\text{full}}, \\ 1 & \text{if } G = G_{\text{diag}} \end{cases} \quad (\text{S.4})$$

and ‘const’ denotes terms that do not depend on  $\mathbf{A}$ . Hence

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

where

$$\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \mathbf{A}} = \left[ \begin{array}{c} -(\xi + 2 - 2\omega_2)/2 \\ -\frac{1}{2} \mathbf{D}_d^T \text{vec} \left( E_{\mathbf{q}}(\Sigma^{-1}) \right) \end{array} \right] \quad (\text{S.5})$$

and  $E_{\mathbf{q}}$  denotes expectation with respect to the normalization of

$$m_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \Sigma}(\Sigma) m_{\Sigma \rightarrow \mathbf{p}(\Sigma|\mathbf{A})}(\Sigma).$$

Let  $q(\Sigma)$  denote this normalized density function. Then  $q(\Sigma)$  is an Inverse-G-Wishart distribution with graph  $G_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \Sigma} \in \{G_{\text{full}}, G_{\text{diag}}\}$  and natural parameter vector  $\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma}$ . From Result 6,

$$E_{\mathbf{q}}(\Sigma^{-1}) = \begin{cases} \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma})_1 + \frac{1}{2}(d+1) \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma})_2 \right) \right\}^{-1} & \text{if } G_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \Sigma} = G_{\text{full}}, \\ \left\{ (\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma})_1 + 1 \right\} \left\{ \text{vec}^{-1} \left( \mathbf{D}_d^{+T} (\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma})_2 \right) \right\}^{-1} & \text{if } G_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \Sigma} = G_{\text{diag}}. \end{cases}$$

Noting that the first factor of  $E_{\mathbf{q}}(\Sigma^{-1})$  is  $(\boldsymbol{\eta}_{\mathbf{p}(\Sigma|\mathbf{A}) \leftrightarrow \Sigma})_1 + \omega_2$ , where  $\omega_2$  is given by (S.4), the first update of  $E_{\mathbf{q}}(\Sigma^{-1})$  in Algorithm 2 is justified. Finally, there is the possible need to adjust for the fact that  $m_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \mathbf{A}}(\mathbf{A})$  is proportional to an Inverse G-Wishart density function with  $G = G_{\text{diag}}$ . This is achieved by the conditional step:

$$\text{If } G_{\mathbf{p}(\Sigma|\mathbf{A}) \rightarrow \mathbf{A}} = G_{\text{diag}} \text{ then } E_{\mathbf{q}}(\Sigma^{-1}) \leftarrow \text{diag} \left\{ \text{diagonal} \left( E_{\mathbf{q}}(\Sigma^{-1}) \right) \right\}.$$

### S.3 Illustrative Example Variational Message Passing Details

The variational message passing approach to fitting and approximate inference for statistical models is still quite a new concept. In this section we provide details on the approach for the illustrative example involving the  $t$  response linear mixed model described in Section 8.

#### S.3.1 Data and Hyperparameter Inputs

Let  $\mathbf{y}$  be the vector of responses as defined in (18). Also, let

$$\mathbf{C} = [\mathbf{X} \ \mathbf{Z}]$$

be the full design matrix, where the matrices  $\mathbf{X}$  and  $\mathbf{Z}$  are as defined in (18). The data inputs are  $\mathbf{y}$  and  $\mathbf{C}$ .

The hyperparameter inputs are

$$\sigma_{\beta}, s_{\sigma^2}, \lambda_{\nu}, s_{\Sigma, 1}, \dots, s_{\Sigma, q} > 0.$$

### S.3.2 Factor to Stochastic Node Parameter Initializations

Initialize  $G_{p(A) \rightarrow A}$  and  $\boldsymbol{\eta}_{p(A) \rightarrow A}$  via a call to Algorithm 1 with hyperparameter inputs:

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

Initialize  $G_{p(a) \rightarrow a}$  and  $\boldsymbol{\eta}_{p(a) \rightarrow a}$  via a call to Algorithm 1 with hyperparameter inputs:

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

Note that the initializations of  $G_{p(A) \rightarrow A}$ ,  $\boldsymbol{\eta}_{p(A) \rightarrow A}$ ,  $G_{p(a) \rightarrow a}$  and  $\boldsymbol{\eta}_{p(a) \rightarrow a}$  are part of the prior impositions for  $\Sigma$  and  $\sigma^2$ . These four factor to stochastic node parameters remain constant throughout the variational message passing iterations.

Initialize

$$\boldsymbol{\eta}_{p(v) \rightarrow v} \leftarrow \begin{bmatrix} 0 \\ -\lambda_v \end{bmatrix}.$$

This initialization of  $\boldsymbol{\eta}_{p(v) \rightarrow v}$  corresponds to the prior imposition for  $v$ . This factor to stochastic node natural parameter remains constant throughout the variational message passing iterations.

The remaining factor to stochastic node natural parameters in the Figure 5 factor graph are updated in the variational message passing iterations, but require initial values. In theory, they can be set to any legal value according to the relevant exponential family. The following initializations, which are used in the code that produced Figure 6, are simple legal natural parameter vectors:

$$G_{p(\Sigma|A) \rightarrow A} \leftarrow G_{\text{diag}}, \quad \boldsymbol{\eta}_{p(\Sigma|A) \rightarrow A} \leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \mathbf{D}_q^T \text{vec}(\mathbf{I}_q) \end{bmatrix},$$

$$G_{p(\Sigma|A) \rightarrow \Sigma} \leftarrow G_{\text{full}}, \quad \boldsymbol{\eta}_{p(\Sigma|A) \rightarrow \Sigma} \leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \mathbf{D}_q^T \text{vec}(\mathbf{I}_q) \end{bmatrix},$$

$$G_{p(\sigma^2|a) \rightarrow a} \leftarrow G_{\text{diag}}, \quad \boldsymbol{\eta}_{p(\sigma^2|a) \rightarrow a} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix},$$

$$G_{p(\sigma^2|a) \rightarrow \sigma^2} \leftarrow G_{\text{full}}, \quad \boldsymbol{\eta}_{p(\sigma^2|a) \rightarrow \sigma^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix},$$

$$\boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma} \leftarrow \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \text{vec}(\mathbf{I}_q) \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})} \leftarrow \begin{bmatrix} \mathbf{0}_{p+mq} \\ -\frac{1}{2} \text{vec}(\mathbf{I}_{p+mq}) \end{bmatrix},$$

$$\boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})} \leftarrow \begin{bmatrix} \mathbf{0}_{p+mq} \\ -\frac{1}{2} \text{vec}(\mathbf{I}_{p+mq}) \end{bmatrix}, \quad \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2} \leftarrow \begin{bmatrix} -2 \\ -1 \end{bmatrix}$$

and

$$\boldsymbol{\eta}_{p(\mathbf{b}|v) \rightarrow v} \leftarrow \begin{bmatrix} 1 \\ -1.1 \end{bmatrix}.$$

The messages involving the  $b_{\ell}$ ,  $1 \leq \ell \leq N$ , nodes do not need to be included here since these messages are subsumed in the calculations used for the natural parameter updates for the model parameters in Algorithm 2 of McLean & Wand (2019).

### S.3.3 Variational Message Passing Iterations

With all factor to stochastic node initialisations accomplished, now we describe the iterative updates inside the variational message passing cycle loop. Each iteration involves:

- updating the stochastic node to factor message parameters.
- updating the factor to stochastic node message parameters.

#### S.3.3.1 Stochastic Node to Factor Message Parameter Updates

The stochastic node to factor message updates are quite simple and follow from, e.g., equation (7) of Wand (2017). For the Figure 5 factor graph the updates are:

$$\begin{aligned}
G_{A \rightarrow p(\Sigma|A)} &\leftarrow G_{p(A) \rightarrow A}, & \boldsymbol{\eta}_{A \rightarrow p(\Sigma|A)} &\leftarrow \boldsymbol{\eta}_{p(A) \rightarrow A}, \\
G_{\Sigma \rightarrow p(\Sigma|A)} &\leftarrow G_{\text{full}}, & \boldsymbol{\eta}_{\Sigma \rightarrow p(\Sigma|A)} &\leftarrow \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma}, \\
\boldsymbol{\eta}_{\Sigma \rightarrow p(\beta, \mathbf{u}|\Sigma)} &\leftarrow \boldsymbol{\eta}_{p(\Sigma|A) \rightarrow \Sigma}, & \boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow p(\beta, \mathbf{u}|\Sigma)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})}, \\
\boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b})} &\leftarrow \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})}, & G_{a \rightarrow p(\sigma^2|a)} &\leftarrow G_{p(a) \rightarrow a}, \\
\boldsymbol{\eta}_{a \rightarrow p(\sigma^2|a)} &\leftarrow \boldsymbol{\eta}_{p(a) \rightarrow a}, & G_{\sigma^2 \rightarrow p(\sigma^2|a)} &\leftarrow G_{\text{full}}, \\
\boldsymbol{\eta}_{\sigma^2 \rightarrow p(\sigma^2|a)} &\leftarrow \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2}, & \boldsymbol{\eta}_{\sigma^2 \rightarrow p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b})} &\leftarrow \boldsymbol{\eta}_{p(\sigma^2|a) \rightarrow \sigma^2}
\end{aligned}$$

and

$$\boldsymbol{\eta}_{v \rightarrow p(\mathbf{b}|v)} \leftarrow \boldsymbol{\eta}_{p(v) \rightarrow v}.$$

Some additional remarks concerning stochastic node to factor updates are:

- The stochastic node to factor messages corresponding to the extremities of the Figure 5 factor graph, such as the message from  $A$  to  $p(A)$ , are not required in the variational message passing iterations. Therefore, updates for these messages can be omitted.
- Some of the stochastic node to factor message parameter updates, such as that for  $\boldsymbol{\eta}_{A \rightarrow p(\Sigma|A)}$ , remain constant throughout the iterations. However, for simplicity of exposition, we list all of the updates together.

#### S.3.3.2 Factor to Stochastic Node Message Parameter Updates

The updates for the parameters of factor to stochastic node messages are a good deal more complicated than the reverse messages. For the illustrative example, these updates are encapsulated in three algorithms across three different articles. Algorithm 2 plays an important role for the variance and covariance matrix parameter parts of the factor graph.

Use Algorithm 2 with:

**Shape Parameter Input:** 1.

**Graph Inputs:**  $G_{\sigma^2 \rightarrow p(\sigma^2|a)}, G_{a \rightarrow p(\sigma^2|a)}$ .

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

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

Use Algorithm 2 with:

**Shape Parameter Input:**  $2q$

**Graph Inputs:**  $G_{\Sigma \rightarrow p(\Sigma|A)}, G_{A \rightarrow p(\Sigma|A)}$

**Natural Parameter Inputs:**  $\boldsymbol{\eta}_{\Sigma \rightarrow \text{p}(\Sigma|A)}$ ,  $\boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow \Sigma}$ ,  $\boldsymbol{\eta}_{A \rightarrow \text{p}(\Sigma|A)}$ ,  $\boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow A}$   
**Outputs:**  $G_{\text{p}(\Sigma|A) \rightarrow \Sigma}$ ,  $\boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow \Sigma}$ ,  $G_{\text{p}(\Sigma|A) \rightarrow A}$ ,  $\boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow A}$

Use the Gaussian Penalisation Fragment of Wand (2017, Section 4.1.4):

**Hyperparameter Input:**  $\sigma_{\beta}^2$

**Natural Parameter Inputs:**  $\boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow \text{p}(\beta, \mathbf{u}|\Sigma)}$ ,  $\boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})}$ ,  
 $\boldsymbol{\eta}_{\Sigma \rightarrow \text{p}(\beta, \mathbf{u}|\Sigma)}$ ,  $\boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma}$

**Outputs:**  $\boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})}$ ,  $\boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma}$

Use the  $t$  Likelihood Fragment of McLean & Wand (2019, Algorithm 2):

**Data Inputs:**  $\mathbf{y}, \mathbf{C}$

**Natural Parameter Inputs:**  $\boldsymbol{\eta}_{(\beta, \mathbf{u}) \rightarrow \text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b})}$ ,  $\boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})}$ ,  
 $\boldsymbol{\eta}_{\sigma^2 \rightarrow \text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b})}$ ,  $\boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2}$ ,  
 $\boldsymbol{\eta}_v \rightarrow \text{p}(\mathbf{b}|v)$ ,  $\boldsymbol{\eta}_{\text{p}(\mathbf{b}|v) \rightarrow v}$

**Outputs:**  $\boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})}$ ,  $\boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2}$ ,  $\boldsymbol{\eta}_{\text{p}(\mathbf{b}|v) \rightarrow v}$

Regarding, the last two fragment updates it should be noted that Wand (2017) and McLean & Wand (2019) work with the “vec” versions of Multivariate Normal and Inverse G-Wishart natural parameter vectors. To match the “vech” natural parameter forms used in Algorithms 1 and 2 of the current article conversions given by (S.1) and (S.2) are required.

### S.3.4 Determination of Posterior Density Function Approximations

After convergence of the variational message passing iterations, the optimal  $q^*$ -densities for each stochastic node are obtained by multiplying each of the messages that pass messages to that node. See, for example, (10) of Wand (2017). We now give details for the model parameters  $\Sigma$ ,  $\sigma^2$ ,  $(\beta, \mathbf{u})$  and  $v$ .

#### S.3.4.1 Determination of $q^*(\Sigma)$

From (10) of Wand (2017):

$$q^*(\Sigma) \propto \exp \left\{ \left[ \begin{array}{c} \log |\Sigma| \\ \text{vech}(\Sigma^{-1}) \end{array} \right]^T \left( \boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow \Sigma} + \boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma} \right) \right\}.$$

It is apparent that  $q^*(\Sigma)$  is an Inverse Wishart density function with natural parameter vector

$$\boldsymbol{\eta}_{q(\Sigma)} \equiv \boldsymbol{\eta}_{\text{p}(\Sigma|A) \rightarrow \Sigma} + \boldsymbol{\eta}_{\text{p}(\beta, \mathbf{u}|\Sigma) \rightarrow \Sigma}.$$

#### S.3.4.2 Determination of $q^*(\sigma^2)$

Using (10) of Wand (2017):

$$q^*(\sigma^2) \propto \exp \left\{ \left[ \begin{array}{c} \log(\sigma^2) \\ 1/\sigma^2 \end{array} \right]^T \left( \boldsymbol{\eta}_{\text{p}(\sigma^2|a) \rightarrow \sigma^2} + \boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2} \right) \right\}.$$

We see that  $q^*(\sigma^2)$  is an Inverse Chi-Squared density function with natural parameter vector

$$\boldsymbol{\eta}_{q(\sigma^2)} \equiv \boldsymbol{\eta}_{\text{p}(\sigma^2|a) \rightarrow \sigma^2} + \boldsymbol{\eta}_{\text{p}(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow \sigma^2}.$$

### S.3.4.3 Determination of $q^*(\beta, \mathbf{u})$

Another application of (10) of Wand (2017) leads to:

$$q^*(\beta, \mathbf{u}) \propto \exp \left\{ \left[ \begin{array}{c} \beta \\ \mathbf{u} \\ \text{vech} \left( \left[ \begin{array}{c} \beta \\ \mathbf{u} \end{array} \right] \left[ \begin{array}{c} \beta \\ \mathbf{u} \end{array} \right]^T \right) \end{array} \right]^T \left( \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})} + \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})} \right) \right\}.$$

We then have  $q^*(\beta, \mathbf{u})$  having a Multivariate Normal density function with natural parameter vector

$$\boldsymbol{\eta}_{q(\beta, \mathbf{u})} \equiv \boldsymbol{\eta}_{p(\beta, \mathbf{u}|\Sigma) \rightarrow (\beta, \mathbf{u})} + \boldsymbol{\eta}_{p(\mathbf{y}|\beta, \mathbf{u}, \sigma^2, \mathbf{b}) \rightarrow (\beta, \mathbf{u})}.$$

### S.3.4.4 Determination of $q^*(v)$

One last application of (10) of Wand (2017) gives:

$$q^*(v) \propto \exp \left\{ \left[ \begin{array}{c} v \log(v) - \log\{\Gamma(v)\} \\ v \end{array} \right]^T \left( \boldsymbol{\eta}_{p(\mathbf{b}|v) \rightarrow v} + \boldsymbol{\eta}_{p(v) \rightarrow v} \right) \right\}.$$

Therefore,  $q^*(v)$  is a Moon Rock density function with natural parameter vector

$$\boldsymbol{\eta}_{q(v)} \equiv \boldsymbol{\eta}_{p(\mathbf{b}|v) \rightarrow v} + \boldsymbol{\eta}_{p(v) \rightarrow v}.$$

## S.3.5 Conversion from Natural Parameters to Common Parameters

A final set of steps involves conversion of the  $q^*$ -densities to common parameter forms.

### S.3.5.1 Conversion of $q^*(\Sigma)$ to Common Parameter Form

The common parameter form of  $q^*(\Sigma)$  is the Inverse-G-Wishart( $G_{\text{full}}, \xi_{q(\Sigma)}, \Lambda_{q(\Sigma)}$ ) density function where

$$\xi_{q(\Sigma)} = -2(\boldsymbol{\eta}_{q(\Sigma)})_1 - 2 \quad \text{and} \quad \Lambda_{q(\Sigma)} = -2\text{vec}^{-1} \left( \mathbf{D}_q^{+T} (\boldsymbol{\eta}_{q(\Sigma)})_2 \right).$$

Alternatively,  $q^*(\Sigma)$  is the Inverse-Wishart( $\kappa_{q(\Sigma)}, \Lambda_{q(\Sigma)}$ ) density function, as defined by (13), where

$$\kappa_{q(\Sigma)} = \xi_{q(\Sigma)} - q + 1.$$

### S.3.5.2 Conversion of $q^*(\sigma^2)$ to Common Parameter Form

The common parameter form of  $q^*(\sigma^2)$  is the Inverse- $\chi^2$ ( $\delta_{q(\sigma^2)}, \lambda_{q(\sigma^2)}$ ) density function where

$$\delta_{q(\sigma^2)} = -2(\boldsymbol{\eta}_{q(\sigma^2)})_1 - 2 \quad \text{and} \quad \lambda_{q(\sigma^2)} = -2(\boldsymbol{\eta}_{q(\sigma^2)})_2.$$

### S.3.5.3 Conversion of $q^*(\beta, \mathbf{u})$ to Common Parameter Form

The common parameter form of  $q^*(\beta, \mathbf{u})$  is the  $N(\boldsymbol{\mu}_{q(\beta, \mathbf{u})}, \boldsymbol{\Sigma}_{q(\beta, \mathbf{u})})$  density function where

$$\boldsymbol{\mu}_{q(\beta, \mathbf{u})} = -\frac{1}{2} \left\{ \text{vec}^{-1} \left( \mathbf{D}_{p+m_q}^{+T} (\boldsymbol{\eta}_{q(\beta, \mathbf{u})})_2 \right) \right\}^{-1} (\boldsymbol{\eta}_{q(\beta, \mathbf{u})})_1$$

and

$$\boldsymbol{\Sigma}_{q(\beta, \mathbf{u})} = -\frac{1}{2} \left\{ \text{vec}^{-1} \left( \mathbf{D}_{p+m_q}^{+T} (\boldsymbol{\eta}_{q(\beta, \mathbf{u})})_2 \right) \right\}^{-1}.$$

Here  $(\boldsymbol{\eta}_{q(\beta, \mathbf{u})})_1$  denotes the first  $p + m_q$  entries of  $\boldsymbol{\eta}_{q(\beta, \mathbf{u})}$  and  $(\boldsymbol{\eta}_{q(\beta, \mathbf{u})})_2$  denotes the remaining entries of the same vector.

### S.3.5.4 Conversion of $q^*(\nu)$ to Common Parameter Form and Conversion to $q^*(\nu)$

Recall that  $q^*(\nu)$  is a Moon Rock density function. The Moon Rock distribution is not as established as the other distributions appearing in this subsection. Nevertheless, the web-supplement of McLean & Wand (2019) defines a random variable  $x$  to have a Moon Rock distribution with parameters  $\alpha > 0$  and  $\beta > \alpha$ , written  $x \sim \text{Moon-Rock}(\alpha, \beta)$ , if the density function of  $x$  is

$$p(x) = \left[ \int_0^\infty \{t^t/\Gamma(t)\}^\alpha \exp(-\beta t) dt \right]^{-1} \{x^x/\Gamma(x)\}^\alpha \exp(-\beta x), \quad x > 0.$$

Therefore,  $q^*(\nu)$  has a Moon-Rock( $\alpha_{q(\nu)}, \beta_{q(\nu)}$ ) density function where

$$\alpha_{q(\nu)} = (\boldsymbol{\eta}_{q(\nu)})_1 \quad \text{and} \quad \beta_{q(\nu)} = -(\boldsymbol{\eta}_{q(\nu)})_2.$$

Explicitly,

$$q^*(\nu) = \left[ \int_0^\infty \{t^t/\Gamma(t)\}^{\alpha_{q(\nu)}} \exp(-\beta_{q(\nu)} t) dt \right]^{-1} \{v^v/\Gamma(v)\}^{\alpha_{q(\nu)}} \exp(-\beta_{q(\nu)} v), \quad v > 0.$$

Lastly, we note that since  $\nu = 2v$  the  $q^*$ -density function of  $\nu$  is

$$q^*(\nu) = \frac{1}{2} \left[ \int_0^\infty \{t^t/\Gamma(t)\}^{\alpha_{q(\nu)}} \exp(-\beta_{q(\nu)} t) dt \right]^{-1} \\ \times \{(\nu/2)^{\nu/2}/\Gamma(\nu/2)\}^{\alpha_{q(\nu)}} \exp\left(-\frac{1}{2}\beta_{q(\nu)}\nu\right), \quad \nu > 0.$$

## Reference

Novomestky, F. (2012). **matrixcalc**: Collection of functions for matrix calculations. R package. <https://CRAN.R-project.org/package=matrixcalc>