

Objective Bayesian Model Selection in Gaussian Graphical Models

BY CARLOS M. CARVALHO

*Graduate School of Business, University of Chicago,
Chicago, Illinois 60637, U.S.A.
carlos.carvalho@chicagogsb.edu*

AND JAMES G. SCOTT

*Department of Statistical Science, Duke University,
Durham, North Carolina 27708-0251, U.S.A.
james@stat.duke.edu*

SUMMARY

This paper presents a default model-selection procedure for Gaussian graphical models that involves two new developments. First, we develop a default version of the hyper-inverse Wishart prior for restricted covariance matrices, called the hyper-inverse Wishart g -prior, and show how it corresponds to the implied fractional prior for covariance selection using fractional Bayes factors. Second, we apply a class of priors that automatically handles the problem of multiple hypothesis testing implied by covariance selection. We demonstrate our methods on a variety of simulated examples, concluding with a real example analysing covariation in mutual-fund returns. These studies reveal that the combined use of a multiplicity-correction prior on graphs and fractional Bayes factors for computing marginal likelihoods yields better performance than existing Bayesian methods.

Some key words: covariance selection; hyper-inverse Wishart distribution; fractional Bayes factors; Bayesian model selection; multiple hypothesis testing.

1 BAYESIAN GRAPHICAL MODELS

1.1 Introduction

Gaussian graphical models are tools for modelling conditional independence relationships, and they offer many practical advantages in high-dimensional problems. They can make computing more efficient by alleviating the need to handle large matrices; they can yield better predictions by fitting sparser models; and they can aid scientific understanding by breaking down a global model into a collection of local models that are easier to parse.

Yet often the graph itself must be inferred from the data, a process often called covariance selection. Our approach to this problem is Bayesian, meaning that two uncertain quantities must be specified: the prior distribution for Σ under each graph, and a prior distribution over different possible graphs.

The first specification is difficult because there is no common covariance matrix shared by all graphs, but rather an entire collection of covariance matrices $\{\Sigma_G\}$ indexed by all possible graphs. Different graphs imply different numbers of free elements in Σ , and so it is not possible to use an improper prior for each Σ_G as one might do for covariance estimation under a fixed graph, since this would leave the resulting model probabilities defined only up to an arbitrary constant. Instead, we must either elicit a subjective prior for each Σ_G —clearly intractable in high dimensions—or we must choose some default proper prior that is neither too vague nor too precise. Regardless, it is clear that any answer will depend on the priors chosen for the various Σ_G 's. To handle this difficulty, we introduce an objective-Bayesian approach using fractional Bayes factors.

The second task—specifying a prior across different graphs—is deceptively easy, yet there are still pitfalls. This is because graphical model selection poses an implicit problem in multiple hypothesis testing, where a null hypothesis is the exclusion of a single edge from the graph. The seemingly objective choice of assigning all graphs equal prior probability will be shown to flag many false-positive edges, and we develop a class of fully Bayesian edge-selection priors to avoid this problem.

1.2 Notation and priors for constrained covariance matrices

An undirected graph is a pair $G = (V, E)$ with vertex set V and edge set $E = \{(i, j)\}$ for some pairs $(i, j) \in V$. Nodes i and j are adjacent, or neighbours, if $(i, j) \in E$. Complete graphs are those having $(i, j) \in E$ for every $i, j \in V$. Complete subgraphs $C \subset V$ are called cliques; two cliques that overlap in a set S are said to have S as a separator. A partition of a graph G into subgraphs (A, S, B) such that $V = A \cup B$, $S = A \cap B$ is complete, and any path from a node in A to a node in B goes through the separator S is called a decomposition. A sequence of subgraphs that cannot be decomposed further are the prime components of a graph; if every prime component

is complete, the graph is said to be decomposable. A graph on p nodes has $m = p(p - 1)/2$ possible edges. All graphs in this paper are assumed to be decomposable.

A Gaussian graphical model or covariance-selection model (Dempster, 1972) uses such a graphical structure to define a set of pairwise conditional-independence relationships on a p -dimensional normally distributed vector $\mathbf{x} \sim N(0, \Sigma)$. The covariance matrix Σ is restricted by its Markov properties; given $\Omega = \Sigma^{-1}$, elements x_i and x_j of the vector \mathbf{x} are conditionally independent, given their neighbors, if and only if $\Omega_{ij} = 0$. If $G = (V, E)$ is an undirected graph describing the joint distribution of \mathbf{x} , $\Omega_{ij} = 0$ for all pairs $(i, j) \notin E$. The covariance matrix Σ is in $M(G)$, the set of all symmetric positive-definite matrices having elements in Σ^{-1} set to zero for all $(i, j) \notin E$.

The hyper-inverse Wishart distribution is a general class of hyper-Markov laws introduced by Dawid & Lauritzen (1993) for a covariance matrix $\Sigma \in M(G)$, where $G = (V, E)$ is a decomposable graph. The notation is $(\Sigma | G) \sim \text{HIW}_G(b, D)$, where $b \in \mathbb{R}^+$ is a degrees-of-freedom parameter, and where $D \in M(G)$ is a symmetric positive-definite scale matrix. The density of this distribution is defined with respect to the product of Lebesgue measures for the (i, j) elements of Σ for which $(i, j) \in E$, subject to the conditions that Σ_C is symmetric and positive definite for all cliques in the junction-tree representation of G . The density of Σ can be obtained from the clique-specific marginal densities as a ratio of products over cliques and separators:

$$p(\Sigma | G) = \frac{\prod_{P \in \mathcal{P}} p(\Sigma_P | b, D_P)}{\prod_{S \in \mathcal{S}} p(\Sigma_S | b, D_S)}, \quad (1)$$

where, for each clique C , $\Sigma_C \sim \text{IW}(b, D_C)$ with density

$$p(\Sigma_C | b, D_C) \propto |\Sigma_C|^{-(b/2+|C|)} \exp \left\{ -\frac{1}{2} \text{tr} (\Sigma_C^{-1} D_C) \right\}. \quad (2)$$

The factorization in (1) holds if we assume that if $S = P_1 \cap P_2$ the elements of Σ_S are common in Σ_{P_1} and Σ_{P_2} . For further details and explanations, refer to Dawid & Lauritzen (1993), Giudici & Green (1999) and Letac & Massam (2007).

1.3 Marginal likelihoods for graphs

Suppose we observe a set of p -dimensional vectors $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, where each $\mathbf{x}_i \sim N(0, \Sigma)$. Let X be the matrix of observations concatenated by row; let X_j refer to column j of X and \mathbf{x}_i to row i ; let X_C refer to the columns of X corresponding to the nodes in clique C ; and assume that $\Sigma \in M(G)$ for some unknown decomposable graph G on p nodes. The posterior probability of a graph G is

$$p(G | X) \propto p(G) \int_{\Sigma \in M(G)} p(X | \Sigma, G) \cdot p(\Sigma | G) d\Sigma, \quad (3)$$

where $p(G)$ is the prior probability of the graph, and where the integral is the marginal likelihood of the data under G . If $\Sigma \sim \text{HIW}_G(b, D)$, the integral in (3) is available in closed form using the ratio of the prior and posterior normalizing constants:

$$p(X | G) = (2\pi)^{-np/2} \frac{h(G, b, D)}{h(G, b^*, D^*)}, \quad (4)$$

where $b^* = b + n$ and $D^* = D + X'X$. The normalizing constant $h(\cdot)$ is

$$h(G, b, D) = \frac{\prod_{P \in \mathcal{P}} \left| \frac{1}{2} D_P \right|^{\frac{(b+|P|-1)}{2}} \Gamma_{|P|} \left(\frac{b+|P|-1}{2} \right)^{-1}}{\prod_{S \in \mathcal{S}} \left| \frac{1}{2} D_S \right|^{\frac{(b+|S|-1)}{2}} \Gamma_{|S|} \left(\frac{b+|S|-1}{2} \right)^{-1}}, \quad (5)$$

where $\Gamma_p(x) = \pi^{p(p-1)/4} \cdot \prod_{j=1}^p \Gamma(x + (1-j)/2)$ is the multivariate gamma function.

2 MODEL-SELECTION PRIORS FOR RESTRICTED COVARIANCE MATRICES

2.1 Criteria for model-selection priors

The expression for the marginal likelihood in (3) involves an integral over the prior for Σ under the graph G . This integral will typically be very sensitive to different choices of the prior, which is a general phenomenon in model-selection problems (Berger & Pericchi, 2001). Unlike in estimation problems, this sensitivity does not diminish as more data are collected, making a smart choice of $p(\Sigma | G)$ for each decomposable graph G the main difficulty in graphical-model selection. In particular, neither improper priors nor vague proper priors can be used.

In all but the smallest of problems, $p(\Sigma | G)$ must be a conjugate hyper-inverse Wishart prior; otherwise it will not be possible to make use of (4) for computing marginal likelihoods. Other priors will require approximating the integrals in (3), and in such a large model space, the need to do so repeatedly will usually pose an insurmountable obstacle. We accept that this practical requirement tethers us to a very restricted class of models, and we seek priors that are as well-behaved as possible under this constraint.

“Well-behaved” is somewhat difficult to judge, since distributions over the space of constrained covariance matrices are not very accessible by intuition. But since every graph implies a complete set of univariate conditional-independence relationships, graphical-model selection can be viewed as simultaneously performing variable selection for all of these conditionals. It is typically easier to assess priors for graphically constrained covariance matrices by studying the properties of the priors they induce on all implied conditional regression models. This point will be developed in §4.

2.2 A conventional proper prior

The most popular choice of prior for use in covariance selection is $\Sigma \sim \text{HIW}_G(\delta, \tau I)$, where the scale matrix $D = \tau I$ is proportional to the identity matrix. We call this “conventional proper prior,” following Berger & Pericchi (2001). Examples of these or similar priors being used to compute marginal likelihoods can be found in Giudici (1996), Giudici & Green (1999), Dobra et al. (2004), Jones et al. (2005), Atay-Kayis & Massam (2005), and Carvalho & West (2007), among others. The scale parameter τ must be chosen to match the expected scale of the data; if τ is too large, the prior for Σ under each graph will wash out the likelihood, and there will be no basis for discriminating among competing graphs.

The now-standard notation of the conventional proper prior can be confusing: a single scale matrix D may be used to specify $p(\Sigma \mid G)$ for all graphs, but this scale matrix means different things under each graph, since different graphs imply different configurations of free elements in Σ . Under G , only the (i, j) elements of D for which the edge $(i, j) \in G$ are relevant in determining the distribution of Σ ; other, nonfree, elements of Σ must be filled in using the (deterministic) completion operation described in Massam & Neher (1998); see also Dawid & Lauritzen (1993) and Carvalho et al. (2007). Hence the notation $\Sigma \sim \text{HIW}_G(b, D)$ must be taken as convenient shorthand for the statement that Σ depends upon the free elements of D implied by the graph G .

2.3 The hyper-inverse Wishart g -prior

Given the practical need for conjugacy, we suggest another possible form of the hyper-inverse Wishart distribution, one where D involves the cross-product matrix:

$$(\Sigma \mid G) \sim \text{HIW}_G(\delta, gX'X),$$

where g is some suitably small fraction such as $1/n$.

We call this the hyper-inverse Wishart g -prior by analogy with Zellner’s g -prior in linear regression (Zellner, 1986), and we recommend it as an alternative to the conventional $\text{HIW}(\delta, \tau I)$ for use in graphical model-selection. The similarity to the g -prior in regression is more than superficial, since this prior will be shown to induce a g -like prior for the univariate conditional regression models implied by G . This along with other theoretical and methodological properties will be examined in §4 and §5.

3 FRACTIONAL BAYES FACTORS FOR COVARIANCE SELECTION

3.1 Motivation for the hyper-inverse Wishart g -prior

Direct use of the hyper-inverse Wishart g -prior for covariance selection is incoherent, since it involves a double use of the data. We now show, however, that the prior

arises very naturally through the use of fractional Bayes factors, and that reusing the data can be avoided.

Fractional Bayes factors were proposed by O’Hagan (1995) as a default Bayesian model-selection technique for use when prior information is weak. The idea is to train a noninformative prior for each model using a small fractional power g of the likelihood function. This is done simultaneously for all models being considered, converting all noninformative priors into proper priors that are then used to select a model with the remainder of the likelihood.

Choose $g \in (0, 1)$ and let $p_N(\Sigma \mid G)$ be a noninformative (typically improper) prior for Σ under a decomposable graph G . The fractional Bayes factor for graphs G_1 and G_2 is then $\text{FBF}_g(G_1, G_2) = Q_g(X \mid G_1)/Q_g(X \mid G_2)$, where

$$Q_g(X \mid G) = \frac{\int p_N(\Sigma \mid G)p(X \mid \Sigma, G) \, d\Sigma}{\int p_N(\Sigma \mid G)p(X \mid \Sigma, G)^g \, d\Sigma}. \quad (6)$$

Then $p^*(\Sigma \mid G) \propto p_N(\Sigma \mid G) \cdot p(X \mid \Sigma, G)^g$ is called the *implied fractional prior*, where the constant of proportionality is the integral of the given expression, and $p(X \mid \Sigma, G)^{1-g}$ is called the *implied fractional likelihood*.

Equation (6) clearly depends upon the choice of a noninformative prior for Σ . The obvious choice given our need for conjugacy is to simply define, for $\Sigma \in M(G)$,

$$p_N(\Sigma \mid G) \propto \frac{\prod_{P \in \mathcal{P}} |\Sigma_P|^{-|P|}}{\prod_{S \in \mathcal{S}} |\Sigma_S|^{-|S|}}, \quad (7)$$

an improper prior that makes use of the same factorization over cliques and separators, and is defined with respect the same measure, as the hyper-inverse Wishart distribution. Interestingly, this prior has clique marginals that correspond to one of many versions of Fisher’s fiducial prior (Sun & Berger, 2007), so called because it yields Fisher’s fiducial distribution for marginal variances. These clique marginals, $p(\Sigma_C) \propto |\Sigma_C|^{-|C|}$, also have the same form as the priors that yield exact frequentist matching for means and variances (Geisser & Cornfield, 1963) when used for covariance estimation. We make no such claims regarding frequentist matching for the prior in (7), and use it only because it is a usefully conjugate generalization of a familiar prior to the space of covariance matrices with graphical structure.

We can now state the main result of this section:

Theorem 3.1. *The hyper-inverse Wishart g -prior $(\Sigma \mid G) \sim \text{HIW}_G(gn, gX'X)$ is the implied fractional prior for Σ corresponding to the prior in (7), where $0 < g < 1$ is the fraction of the likelihood used for training. The fractional marginal likelihood is*

$$Q_g(X \mid G) = (2\pi)^{-np/2} \frac{h(G, gn, gX'X)}{h(G, n, X'X)}, \quad (8)$$

with $h(G, b, D)$ defined as in (5).

Proof. Follows immediately from the conjugacy of the hyper-inverse Wishart prior with the normal likelihood, and from Equation (4). \square

We emphasize that this procedure does not merely specify a single prior distribution, but rather a whole cohort of objective prior distributions for all $\Sigma_G \in M(G)$.

3.2 Choice of g

An obvious issue with the use of this methodology is the choice of g . If the hyper-inverse Wishart g -prior were used as a real prior, it would be possible to place a hyperprior on g , and not tie ourselves down to a specific value; indeed, Liang et al. (2008) recommend exactly this approach toward g -priors in linear-model selection.

Yet when interpreting the hyper-inverse Wishart g -prior in terms of fractional Bayes factors, it is no longer possible to put a prior on g . This is because g is not a model parameter about which there is information in the likelihood, but rather the fractional power of the likelihood itself used for training the noninformative prior in (7). This fraction must be chosen outright in order for the fractional marginal likelihoods in (8) to be well-defined.

Several criteria help to guide this choice. First, there is an established tradition of using “minimal training sample” sizes to calibrate default Bayes factors; see, for example, O’Hagan (1995), Berger & Pericchi (1996), and Berger & Pericchi (2001). A minimal training sample is the smallest sample size needed to convert all improper priors such as (7) into proper priors. The intuition is that as much of the data as possible should be held back to choose between models. It is easy to see from (1) and (2) that the minimal training sample size is 1, suggesting that g be $1/n$.

Second, it is clear that g must be $O(1/n)$ in order for the implied fractional prior to correspond asymptotically to the gold standard of a carefully elicited subjective prior distribution. If g decreases too slowly as a function of n , the implied fractional prior will asymptotically overwhelm the likelihood; if it decreases too fast, the prior will become arbitrarily diffuse. Neither behavior could possibly result from the choices of a careful elicitee making intelligent decisions about each $p(\Sigma | G)$. (We find the hypothetical “careful elicitation” a useful ideal to keep in mind, even if dimensionality makes this ideal impossible to attain.)

Finally, the implied fractional prior for Σ should have heavy tails. Choosing $gn = 1$ implies that the vector \mathbf{x} is marginally Cauchy, and that each $p(\Sigma | G)$ is heavy-tailed without being too vague. This choice dovetails with the advice given by Liang et al. (2008), who themselves generalize the recommendations of Jeffreys (1961) and Zellner & Siow (1980).

As a default choice, we recommend setting $g = 1/n$, though this is not a hard rule, and other choices that decay like $1/n$ may be reasonable (and can be judged by the reasonableness of their effect on the implied fractional prior). Robustness to these choices should be considered, just as it should be in subjective analyses.

4 PROPERTIES OF THE HYPER-INVERSE WISHART g -PRIOR

4.1 Information consistency

The consistency of fractional Bayes factors as $n \rightarrow \infty$ is a well-known result from O'Hagan (1995). We instead consider a second notion of consistency, often called *information consistency* or *finite-sample consistency*, that describes how a Bayes factor behaves for fixed n with respect to a test statistic that would be used to perform a classical test of significance on the same problem.

The canonical example of an information-inconsistent procedure is model selection in linear regression using fixed- g versions of Zellner's g -prior (Zellner & Siow, 1980; Liang et al., 2008). Imagine testing a specific regression model M_A having k possible covariates against the null model M_0 having only an intercept term. If the usual F statistic for testing M_A against M_0 goes to infinity for fixed n and $k < n - 1$, the evidence against M_0 is overwhelming, and one would expect the Bayes factor $\text{BF}(M_A : M_0)$ to diverge.

But under the standard g -prior, this Bayes factor instead converges to the fixed constant $(1 + g)^{(n-k-1)/2}$. This gives an intrinsic limitation, one that is not shared by the F statistic, to how strongly the Bayes factor may support the bigger model. Such behavior is intuitively unappealing, since $\text{pr}(F > C \mid M_0) \rightarrow 0$ as $C \rightarrow \infty$.

A natural question is whether the fractional Bayes factors defined above exhibit a similar information paradox. We show that they do not, in two related senses.

4.2 Tests against the null graph

Let G_0 denote the null graph having no edges, and let G_A denote the graph to be compared with the null. The Bayes factor for comparing these two models is

$$\begin{aligned} \text{BF}(G_0 : G_A) = K \cdot & \frac{\prod_{j=1}^p \left| \frac{g}{2} X'_j X_j \right|^{\frac{gn}{2}}}{\prod_{j=1}^p \left| \frac{1}{2} X'_j X_j \right|^{\frac{n}{2}}} \cdot \frac{\prod_{S \in \mathcal{S}} \left| \frac{g}{2} X'_S X_S \right|^{\frac{gn+|S|-1}{2}}}{\prod_{P \in \mathcal{P}} \left| \frac{g}{2} X'_P X_P \right|^{\frac{gn+|P|-1}{2}}} \\ & \cdot \frac{\prod_{P \in \mathcal{P}} \left| \frac{1}{2} X'_P X_P \right|^{\frac{n+|P|-1}{2}}}{\prod_{S \in \mathcal{S}} \left| \frac{1}{2} X'_S X_S \right|^{\frac{n+|S|-1}{2}}}, \end{aligned} \quad (9)$$

where g is fixed, \mathcal{P} and \mathcal{S} are the prime components and separators of G_A , and the leading term K is

$$K = \left[\frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{gn}{2}\right)} \right]^p \cdot \frac{\prod_{S \in \mathcal{S}} \Gamma_{|S|}\left(\frac{n+|S|-1}{2}\right)}{\prod_{P \in \mathcal{P}} \Gamma_{|P|}\left(\frac{n+|P|-1}{2}\right)} \cdot \frac{\prod_{P \in \mathcal{P}} \Gamma_{|P|}\left(\frac{gn+|P|-1}{2}\right)}{\prod_{S \in \mathcal{S}} \Gamma_{|S|}\left(\frac{gn+|S|-1}{2}\right)}.$$

It remains to define a suitable test statistic b as a basis for assessing information

consistency. Following Lauritzen (1996), let $\widehat{\Omega}_0$ be the maximum-likelihood estimate for the precision matrix under G_0 , and let $\widehat{\Omega}_A$ be the maximum-likelihood estimate under G_A . Then there is a nested sequence $G_0 \subset \dots \subset G_d = G_A$ of decomposable graphs that differ only by a single edge. Let e_i denote the edge in G_i but not in G_{i-1} , and let C_i be the (unique) clique of G_i containing e_i . Then we have the following proposition (see Proposition 5.14 of Lauritzen, 1996, for a proof):

Proposition 4.1. *The test of significance for G_A against the null graph G_0 can be performed by rejecting G_0 for sufficiently small values of $b = |\widehat{\Omega}_0|/|\widehat{\Omega}_A|$. Under G_0 , b is distributed as the product of independent beta random variables $B_1 \cdots B_d$, with $B_i \sim \text{BE}\{(n - |C_i|)/2, 1/2\}$.*

This defines the relevant test statistic b . We now give a precise statement of information consistency for the fractional Bayes factors in Theorem 3.1.

Theorem 4.2. *Let G_A be a decomposable graph having prime components \mathcal{P} , let G_0 be the null graph, and let $\text{FBF}_g(G_A : G_0)$ be the fractional Bayes factor, given data X , corresponding to the noninformative prior in (7). For any finite $n > \max_{P \in \mathcal{P}} |P|$ and for any $0 < g < 1$, $\text{FBF}_g(G_0 : G_A) \rightarrow 0$ as $b \rightarrow 0$.*

Proof. The Bayes factor in (9) simplifies to

$$K \cdot \left(\frac{1}{g}\right)^{(S_1 - gnp)/2} \left(\frac{1}{2}\right)^{(n - gn)(S_2 - p)/2} \cdot \prod_{j=1}^p |X'_j X_j|^{-(n - gn)/2} \cdot \frac{\prod_{P \in \mathcal{P}} |X'_P X_P|^{(n - gn)/2}}{\prod_{S \in \mathcal{S}} |X'_S X_S|^{(n - gn)/2}},$$

where the exponent terms S_1 and S_2 are

$$\begin{aligned} S_1 &= \sum_{P \in \mathcal{P}} |P| \cdot (gn + |P| - 1) - \sum_{S \in \mathcal{S}} |S| \cdot (gn + |S| - 1) \\ S_2 &= \sum_{P \in \mathcal{P}} |P| - \sum_{S \in \mathcal{S}} |S|. \end{aligned}$$

Now apply the formula of Lauritzen (1996) for the determinant of $\widehat{\Omega}_G$, which will exist due to the restriction that $n > \max_{P \in \mathcal{P}} |P|$:

$$|\widehat{\Omega}_G| = n^p \frac{\prod_{S \in \mathcal{S}} |X'_S X_S|}{\prod_{P \in \mathcal{P}} |X'_P X_P|}.$$

This gives

$$\text{BF}(G_0 : G_A) = C \cdot \left(\frac{|\widehat{\Omega}_0|}{|\widehat{\Omega}_A|}\right)^{(n - gn)/2},$$

where C is a fixed, finite term involving g , p , n , and the structure of the graph G_A . The proof of information consistency now follows immediately by plugging the

test statistic b into the above equation, and noticing that for fixed $0 < g < 1$, $\text{BF}(G_0 : G_A) \rightarrow 0$ as $b \rightarrow 0$. \square

4.3 Tests for an implied conditional regression model

Information consistency is important in a second sense: nonzero entries in a precision matrix imply a set of nonzero conditional regression coefficients for each element of \mathbf{x} upon the other elements, and Bayes factors for covariance selection perform variable selection on all of these implied regressions simultaneously. We find that observing the behavior of these implied conditional regression models is a useful window on the behavior of $p(\Sigma | G)$, which is far harder to understand intuitively.

The following lemma characterizes these univariate conditionals.

Lemma 4.3. *Suppose $(\mathbf{x} | \Sigma) \sim \text{N}(0, \Sigma)$ and $\Sigma \sim \text{HIW}_G(b, D)$ for some decomposable graph G . Suppose $\mathbf{x} = (z, \mathbf{y})'$ where z is a scalar, and that Σ and D are partitioned*

$$\Sigma = \begin{pmatrix} \Sigma_{zz} & \Sigma_{z\mathbf{y}} \\ \Sigma_{\mathbf{y}z} & \Sigma_{\mathbf{y}\mathbf{y}} \end{pmatrix}, \quad D = \begin{pmatrix} D_{zz} & D_{z\mathbf{y}} \\ D_{\mathbf{y}z} & D_{\mathbf{y}\mathbf{y}} \end{pmatrix}.$$

Then:

- (i) $\Sigma_{z|\mathbf{y}}^{-1} = (\Sigma_{zz} - \Sigma_{z\mathbf{y}}\Sigma_{\mathbf{y}\mathbf{y}}^{-1}\Sigma_{\mathbf{y}z})^{-1} \sim \text{GA}\left(\frac{b+k}{2}, \frac{D_{z|\mathbf{y}}}{2}\right)$
- (ii) $(\Sigma_{z\mathbf{y}}\Sigma_{\mathbf{y}\mathbf{y}}^{-1} | \Sigma_{z|\mathbf{y}}) \sim \text{N}(D_{z\mathbf{y}}D_{\mathbf{y}\mathbf{y}}^{-1}, \Sigma_{z|\mathbf{y}}D_{\mathbf{y}\mathbf{y}}^{-1})$

with k representing the number of neighbors of z under G .

The lemma is proven in Appendix A. We now apply it to give the following theorem, which is intended to be understood as if the implied fractional prior were a true prior, and the implied fractional likelihood were a true likelihood. Since the fractional prior has the exact functional form of a (proper) hyper-inverse Wishart prior, and the fractional likelihood has the functional form of a Gaussian likelihood, there is no mathematical ambiguity in operationally defining “prior” and “likelihood” this way, even if the interpretation is difficult from a pure subjectivist viewpoint.

Theorem 4.4. *Let $X = (\mathbf{z} \ Y)$ be the matrix of observed data having rows $X_i = (z_i, \mathbf{y}_i)$. Let G_A be a decomposable graph having prime components \mathcal{P} , and let M_A be the conditional regression model for z in terms of Y implied by the neighbors of z in G_A , and let M_0 be the null regression model for z . Let F be the usual F -statistic for testing M_A against M_0 , and let $\text{BF}_g(M_A : M_0)$ be the likelihood ratio $p(\mathbf{z}|Y, M_A)/p(\mathbf{z}|Y, M_0)$, where these marginals are defined under the fractional prior and the fractional likelihood in (6) for a fixed g . For any finite $n > \max_{P \in \mathcal{P}} |P|$ and for any $0 < g < 1$, $F \rightarrow \infty \implies \text{BF}_g(M_A : M_0) \rightarrow \infty$.*

Proof. Let Y_z denote the columns of the matrix X corresponding to the neighbors of z under M_A . Applying Lemma 4.3 under the assumption that $n > \max_{P \in \mathcal{P}} |P|$, the hyper-inverse Wishart g -prior gives the following regression relationship:

$$\mathbf{z} = Y_z \mathbf{f} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{N}(0, \phi^{-1}I) \quad (10)$$

$$(\mathbf{f} \mid \phi, M_A) \sim \text{N}\left(\hat{\mathbf{f}}, (g\phi)^{-1}(Y_z'Y_z)^{-1}\right) \quad (11)$$

$$(\phi \mid M_A) \sim \text{Ga}\left(\frac{gn+k}{2}, \frac{gr}{2}\right), \quad (12)$$

where ϕ is the conditional precision or $\Sigma_{\mathbf{z}|Y_z}^{-1}$; I is the $n \times n$ identity matrix, $\hat{\mathbf{f}}' = (Y_z'Y_z)^{-1}Y_z'\mathbf{z}$ is the traditional least-squares estimate for \mathbf{f} ; and $r = \mathbf{z}'(I - P_{Y_z})\mathbf{z}$ with P_{Y_z} denoting the perpendicular projection matrix onto the column space of Y_z . Hence r is the residual sum of squares after regressing \mathbf{z} upon Y_z .

Marginalizing over \mathbf{f} and ϕ , taking care to use the fractional likelihood rather than the full likelihood, gives

$$P(\mathbf{z} \mid Y, M_A) = (\pi)^{-n/2} g^{\frac{gn+2k}{2}} (1-g)^{n/2} \cdot \frac{\Gamma\left(\frac{n+gn+k}{2}\right)}{\Gamma\left(\frac{gn+k}{2}\right)} \cdot r^{-n/2}. \quad (13)$$

Assuming $0 < g < 1$, the relevant Bayes factor can then be computed by recognizing the null model M_0 as a special case of (13) with $k = 0$ and $r = \mathbf{z}'\mathbf{z}$:

$$\text{BF}(M_A : M_0) = C (1 - R_{M_A}^2)^{-n/2}, \quad (14)$$

where C is a fixed term involving g , n , and k , and where $R_{M_A}^2$ is the usual coefficient of determination for model M_A . As the F -statistic grows without bound, $R_{M_A}^2 \rightarrow 1$, and the Bayes factor in (14) clearly diverges. \square

The hierarchical model in (10), (11), and (12) is immediately recognizable as a modified form of Zellner's g -prior for the vector of conditional regression coefficients. Unlike the g -prior, however, this procedure avoids the information paradox.

5 FRACTIONAL MARGINAL LIKELIHOODS: A SIMULATION STUDY

We now study the behavior of the fractional marginal likelihoods through simulations that compare models of differing complexity. The baseline for comparison will be the conventional alternative, the $\text{HIW}_G(\delta, \tau I)$ prior. We hope to illustrate that the casual use of the conventional prior undermines our ability to choose between competing models, and that the objective procedure presented so far is well-behaved.

Data sets of various sizes ($n = 20, 50, 100, 500, 1000$) were simulated from the true model: each $\mathbf{x}_i \sim \text{N}(0, \Sigma)$, where Σ was the 50-dimensional correlation matrix of a stationary Gaussian AR(10) process. This represents a Gaussian graphical model due

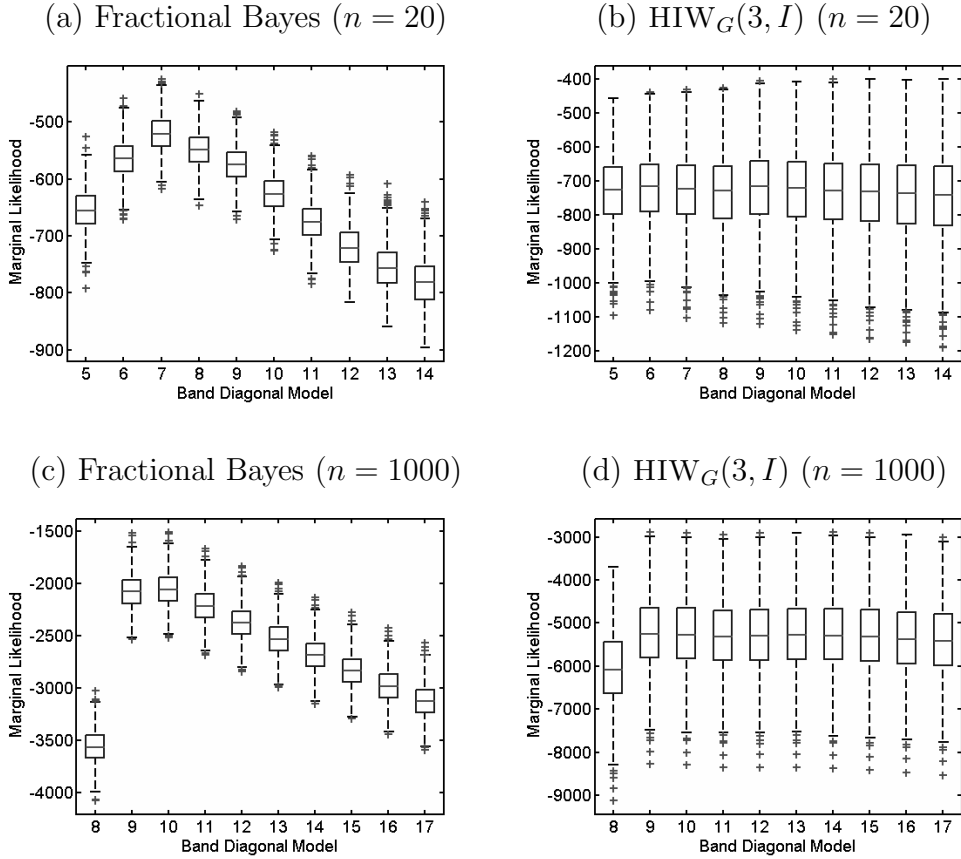


Figure 1: Boxplots of realized marginal likelihoods in the simulation study. The x -axis is the bandwidth of the precision matrix; the y -axis is marginal log-likelihood.

to the band-diagonal form of the precision matrix. Note that $p = 50$ is the number of nodes (i.e. the length of each observation \mathbf{x}_i), but that n is the number of such independent draws observed from the true model.

An appropriate choice for the conventional proper prior is $\Sigma \sim HIW_G(3, I)$, where our choice of $\delta = 3$ reflects the standard advice to give $p(\Sigma | G)$ a finite first moment (Jones et al., 2005). We have chosen $\tau = 1$, since Σ is known to be a correlation matrix. Typically τ is very hard to specify without looking at the data, making this choice, if anything, overly favorable to the conventional prior.

For each sample size, we simulated 1000 datasets from the true model and computed marginal likelihoods for 21 different candidate graphs corresponding to band-diagonal precision matrices of bandwidth 0 through 20 (with the true model having bandwidth 10). Figure 1 gives the frequency distributions of marginal log-likelihoods for each candidate model, which show substantially better separation under the fractional Bayes approach. Results are given for both the smallest ($n = 20$) and largest ($n = 1000$) sample sizes, although the same pattern emerged for all sample sizes.

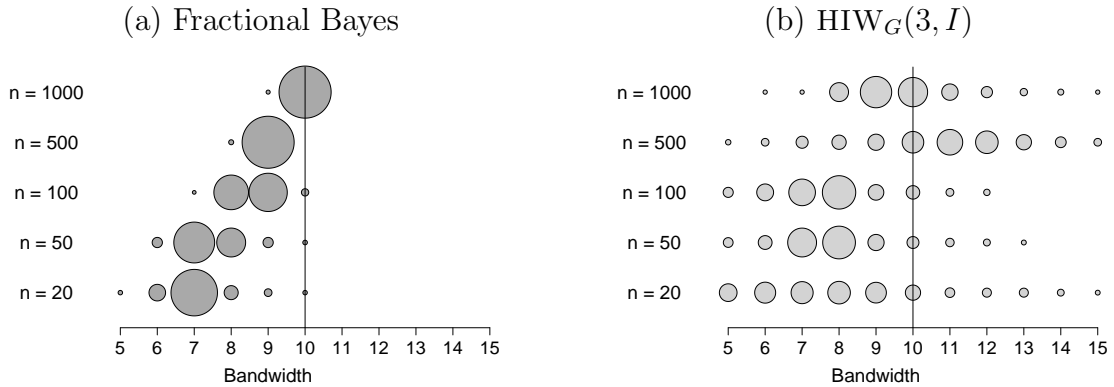


Figure 2: The empirical distribution of the chosen bandwidth under repeated sampling from the true model in the simulation study. The area of each circle represents the fraction out of 1000 independent data sets in which each model had the highest marginal likelihood under the given prior.

For each data set, we also recorded the model with the highest marginal likelihood. This gives a Monte-Carlo estimate for the frequency distribution of the preferred band-diagonal size under the two priors (Figure 2). As n grows, the fractional Bayes factors favor the true model (bandwidth 10) with increasing accuracy, whereas the results from the $HIW_G(3, I)$ are highly erratic.

It is clear that, unlike the conventional prior, our objective procedure prefers sparse models in the absence of enough data to justify extra edges. Yet as the sample size increases, the fractional approach favors more complex models, and eventually chooses the true one almost every time. The conventional prior does not exhibit this tendency nearly as strongly, displaying an unintuitively high level of variation in the choice of model.

One might imagine that the conventional prior, by shrinking the covariance structure toward the identity matrix with its strong pattern of off-diagonal zeros, would yield systematically smaller models. This expectation is not confirmed by our simulation study, indicating that intuitions about shrinkage gleaned from covariance estimation do not necessarily apply to covariance selection. This can be understood in terms of the Ockham’s-razor property of Bayesian marginal likelihoods (Jefferys & Berger, 1992). The conventional prior spreads its mass out quite broadly—an advantage in estimation problems, but crippling in model selection due to the lack of predictions sharp enough to yield any posterior separation of models.

The results from this simulation are consistent with the theoretical results developed so far and strengthen our claim that the fractional approach is indeed an appropriate default procedure for Gaussian-graphical-model selection. We note that the conventional prior induces a set of ridge-regression priors on the complete conditional

models considered in §4.3 (as can be shown through a straightforward application of Lemma 4.3). The problems with ridge-regression priors for variable selection are well understood (Zellner & Siow, 1980; Liang et al., 2008), and give some intuition as to why the conventional $\text{HIW}_G(\delta, \tau I)$ prior is suboptimal for covariance selection.

6 PRIORS OVER GRAPHS

There are two simple approaches to assigning prior probabilities to graphs themselves. The first is to give every graph the same prior probability κ^{-1} , where κ is the number of decomposable graphs on p nodes (which is not trivial to compute). As Giudici & Green (1999) note, this prior is quite heavily concentrated on graphs of middling size due to combinatorial explosion.

The second alternative, rapidly becoming the standard, is to model edge inclusions as independent Bernoulli events with common success probability r (Dobra et al., 2004; Jones et al., 2005). This yields priors of the form

$$p(G) \propto r^k (1 - r)^{m-k}. \quad (15)$$

The constant of proportionality $2^m/\kappa$ is hard to compute, but is the same for all models and can thus be ignored. If the expected fraction of included edges is known quite precisely, this framework may be attractive. Yet often this fraction is not known, making an arbitrary choice of r seem heavy-handed.

Instead, we recommend treating r as a model parameter to be estimated from the data rather than as a fixed tuning constant. This shrinks the graph size to a data-determined value of r , which is estimated from the prevailing edge-inclusion fraction of models with favorable marginal likelihoods.

One possibility is to estimate r by empirical-Bayes methods. George & Foster (2000) give an EM procedure that can be used to compute a maximum-likelihood estimate \hat{r} in linear regression models, but note that it can prove computationally intractable in high dimensions, making it even less suited to graphical model spaces.

A second possibility is to place a prior on r , which turns out to be the easier and more attractive option. Assuming the conjugate beta prior $r \sim \text{BE}(a, b)$ allows an explicit marginalization:

$$p(G) \propto \int_0^1 p(G | r) p(r) dr \propto \frac{\beta(a + k, b + m - k)}{\beta(a, b)}$$

where $\beta(\cdot, \cdot)$ is the beta function. For the default choice of $a = b = 1$, implying a uniform prior on r , this becomes

$$p(G) \propto \frac{(k)!(m-k)!}{(m+1)(m!)} = \frac{1}{m+1} \binom{m}{k}^{-1}. \quad (16)$$

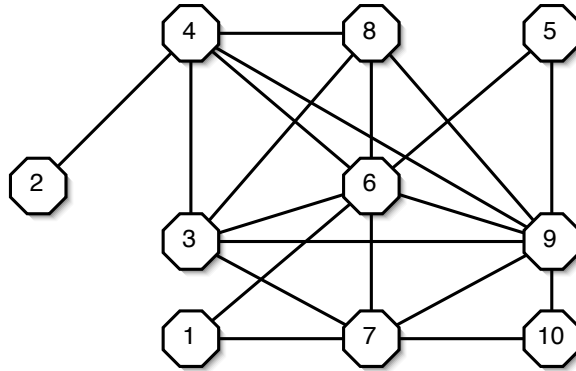


Figure 3: The true 10-node, decomposable graph used in the multiplicity-correction study. All noise nodes were completely unconnected both from this graph and from each other, meaning that any edges involving them are false positives.

Scott & Berger (2006) call these “multiplicity correction” priors in the context of testing exchangeable normal means: if the true k is held fixed while m grows, the expected number of false positives will remain constant. In our context, this will mean an automatic penalty for testing irrelevant edges.

We demonstrate through simulation that the same effect holds here. Beginning with a correlation matrix corresponding to the 10-node graph from Figure 3, we added progressively more “noise” nodes—that is, nodes unconnected both from the true graph and from each other, but that lead to combinatorial explosion in the number of edges that must be tested for inclusion. The numbers of noise nodes chosen were 5, 15, and 40, which in addition to the 10 connected nodes in the true graph imply 60, 300, and 1225 separate hypothesis tests, respectively. In all cases the number of true hypotheses remained fixed at 22, one for each edge in the 10-node graph.

Three sets of tests were performed under each of three different priors on models. These results are summarized in Table 1. Here, “Fully Bayes” uses the model probabilities from (16), while “Oracle Bayes” involves plugging the true value of r into (15) to compute prior graph probabilities. All marginal likelihoods were computed using fractional Bayes factors, with $g = 1/n$.

Notice that the fully Bayesian multiplicity-correction prior in (16) squelches false positives very effectively. This phenomenon is well understood in linear models (Scott & Berger, 2008), but ours is the first demonstration of such an automatic-correction effect in graphical models. The difference between corrected and uncorrected versions is substantial; in the 50-node example, giving all models the same prior probability yields 40 false positives (inclusion probability greater than 50%), whereas imposing the multiplicity-correction prior yields none. Importantly, this approach does not depend upon the choice of an arbitrary hyperparameter r , although if subjective inputs are required, they can be accommodated through a different beta prior while

Edge	Number of Noise Nodes								
	No Correction			Oracle Bayes			Fully Bayes		
	5	15	40	5	15	40	5	15	40
(1,6)	99	99	99	99	99	99	99	99	99
(3,4)	16	2	0	1	0	3	1	1	1
(3,6)	99	99	99	99	99	97	99	99	99
(3,8)	18	0	0	10	2	0	6	1	0
(3,9)	31	3	0	0	0	0	0	0	0
(4,6)	99	99	99	99	99	99	99	99	99
(4,9)	99	99	99	99	99	99	99	99	99
(5,6)	16	5	0	22	21	18	24	23	22
(5,9)	36	14	31	31	30	34	31	31	33
(6,7)	58	76	74	14	2	0	8	3	0
(6,9)	22	3	0	1	0	0	1	0	0
(8,9)	43	4	7	14	2	0	7	1	0
(9,10)	89	95	99	71	69	76	60	60	69
FPs:	6	11	40	1	1	1	1	1	0

Table 1: Estimated inclusion probabilities for specific edges as the number of unconnected noise nodes grows. The last line of the table shows the number of falsely positive flags—that is, $\geq 50\%$ edge inclusion probability—among other, non-enumerated edges. All probabilities were calculated using 5 million iterations of the FINCS algorithm (Scott & Carvalho, 2008).

still retaining closed-form answers.

Other differences from standard priors also emerge. Consider, for example, edges (5, 9) and (6, 7). If all models are given equal prior probabilities, adding more noise edges makes (6, 7) appear stronger and (5, 9) appear slightly weaker. Yet the opposite happens using the multiplicity-correction priors: the addition of more noise nodes makes (6, 7) disappear entirely and yet retains (5, 9) at close to its original strength.

This behavior suggests that (16) does not merely shrink edge-inclusion probabilities to 0 uniformly as k remains fixed and m grows. Rather, it differentially rewards edges that participate in more parsimonious models, suggesting a fundamental difference from (15) in the way mass is apportioned across model space.

7 EXAMPLE: MUTUAL FUNDS

A crucial input for many dynamic portfolio-selection problems is the estimated covariance matrix Σ for a collection of asset returns. Naive procedures can often yield unstable estimates of Σ , but as Carvalho & West (2007) show, graphical models offer a potent tool for shrinkage.

We now illustrate the proposed methodology on a set of 86 monthly returns for 59 mutual funds from a variety of different sectors: 13 U.S. bond funds, 30 U.S. stock funds, 7 balanced funds investing in both U.S. stocks and bonds, and 9 international

Method		Predictive Squared Error	
Likelihood	Multiplicity	Top Model	Model Average
HIW- g	Sector-specific	554.654	542.342
HIW- g	Global	623.550	611.214
HIW- g	No Correction	642.874	636.176
HIW($3, \tau I$)	No Correction	679.643	661.077
Saturated Model		1,749.072	

Table 2: Sum of squared errors in predicting missing data in the mutual-fund example. All model searches used 3 million iterations of FINCS on the 60-month training set. Model-averaged predictions are averaged over the top 500 models discovered during the course of the search.

stock funds. The monthly returns were split into a 60-month training set and a 26-month prediction set. We study the predictive performance of each of the different graphical-modelling regimes considered here, incorporating different combinations of priors for Σ and priors over graphs.

The last of these regimes involved expanding the prior specification in (16) to include several different parameters for the prior inclusion probabilities for edges in each of the different fund sectors. This reflects our prior understanding that different patterns of covariation might prevail in different asset classes, and that it may be suboptimal to force edges participating in these different covariational patterns to shrink towards some common inclusion probability. Instead, we allow block-by-block patterns of (possibly differing) sparsity to emerge. Note that global shrinkage could still be accomplished via a hierarchical prior; here we have simply chosen to give each inclusion probability an independent uniform prior.

For each prior specification, we used the FINCS algorithm (Scott & Carvalho, 2008) on the 60-month training set to search for good models and compute posterior means $\hat{\Sigma}$. These means were then used to predict observations in the 26-month validation set, where in each month the 56 “observed” returns predicted the remaining 3 “missing” returns (which are known in reality). We repeated this process for all of the $\binom{59}{3} = 32,509$ combinations of 3 unobserved columns, and then computed the total squared-error in imputing the missing values.

When all models were given equal prior probability, the fractional approach yielded a 5.5% improvement in total mean-squared error (642.9 versus 679.6) under the top models discovered (see Table 2). Applying the multiplicity-correction priors in (16) reduced the squared-error to 623.6, while the block-by-block prior (with different r parameters corresponding to different asset classes) yielded a further reduction to 554.7. These may seem like small improvements, but are actually quite substantial given the highly unpredictable nature of financial markets.

These results demonstrate that each of the methodological advancements con-

sidered here can improve predictive accuracy in covariance selection problems. It is especially interesting to note the strength of the sector-specific multiplicity correction, suggesting that the ease with which our prior specification accommodates structural information can be a real advantage in practice.

8 DISCUSSION

We have introduced a new method of covariance selection based upon objective Bayesian ideas. The strengths of our approach are the theoretical guarantees of §4, the intuitive behavior of fractional marginal likelihoods demonstrated in §5, and the strong control over false positives shown in §6. Our prior specification also accommodates subjective information about the structure of the problem in a flexible and intuitive way, which may, as in the case of the mutual-fund example, improve predictive accuracy. Moreover, these theoretical developments ensure that such predictions rest upon the firm ground of a well-behaved set of posterior model probabilities.

REFERENCES

- ATAY-KAYIS, A. & MASSAM, H. (2005). A Monte Carlo method for computing the marginal likelihood in nondecomposable Gaussian graphical models. *Biometrika* **92**, 317–35.
- BERGER, J. O. & PERICCHI, L. (1996). The intrinsic Bayes factor for model selection and prediction. *J. Amer. Statist. Assoc.* **91**, 109–122.
- BERGER, J. O. & PERICCHI, L. (2001). Objective Bayesian methods for model selection: introduction and comparison. In *Model Selection*, vol. 38 of *Institute of Mathematical Statistics Lecture Notes – Monograph Series*. Beachwood, pp. 135–207.
- CARVALHO, C., MASSAM, H. & WEST, M. (2007). Simulation of hyper-inverse Wishart distributions in graphical models. *Biometrika* **94**, 647–59.
- CARVALHO, C. & WEST, M. (2007). Dynamic matrix-variate graphical models. *Bayesian Anal.* **2**, 69–96.
- COWELL, R. G., DAWID, A. P., LAURITZEN, S. L. & SPIEGELHALTER, D. J. (1999). *Probabilistic Networks and Expert Systems*. Springer-Verlag New York, 1st Ed.
- DAWID, A. P. & LAURITZEN, S. L. (1993). Hyper-Markov laws in the statistical analysis of decomposable graphical models. *Ann. Statist.* **21**, 1272–317.
- DEMPSTER, A. (1972). Covariance selection. *Biometrics* **28**, 157–75.

- DOBRA, A., JONES, B., HANS, C., NEVINS, J. & WEST, M. (2004). Sparse graphical models for exploring gene expression data. *J. Mult. Anal.* **90**, 196–212.
- GEISSER, S. & CORNFIELD, J. (1963). Posterior distributions for multivariate normal parameters. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **25**, 368–76.
- GEORGE, E. I. & FOSTER, D. P. (2000). Calibration and empirical Bayes variable selection. *Biometrika* **87**, 731–47.
- GIUDICI, P. (1996). Learning in graphical Gaussian models. In *Bayesian Statistics 5*, Eds. J. Berger, J. Bernardo, A. Dawid & A. Smith. Oxford University Press, pp. 621–8.
- GIUDICI, P. & GREEN, P. J. (1999). Decomposable graphical Gaussian model determination. *Biometrika* **86**, 785–801.
- JEFFERYS, W. & BERGER, J. (1992). Ockham’s razor and Bayesian analysis. *Am. Sci.* **80**, 64–72.
- JEFFREYS, H. (1961). *Theory of Probability*. Oxford University Press, 3rd Ed.
- JONES, B., CARVALHO, C., DOBRA, A., HANS, C., CARTER, C. & WEST, M. (2005). Experiments in stochastic computation for high-dimensional graphical models. *Statist. Sci.* **20**, 388–400.
- LAURITZEN, S. L. (1996). *Graphical Models*. Oxford: Clarendon Press.
- LETAC, G. & MASSAM, H. (2007). Wishart distributions on decomposable graphs. *Ann. Statist.* **35**, 1278–1323.
- LIANG, F., PAULO, R., MOLINA, G., CLYDE, M. & BERGER, J. (2008). Mixtures of g -priors for Bayesian variable selection. *J. Amer. Statist. Assoc.* **103**, 410–23.
- MASSAM, H. & NEHER, E. (1998). Estimation and testing for lattice conditional independence models on euclidean jordan algebras. *Ann. Statist.* **26**, 1051–82.
- O’HAGAN, A. (1995). Fractional Bayes factors for model comparison. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **57**, 99–138.
- ROVERATO, A. (2000). Cholesky decomposition of a hyper-inverse Wishart matrix. *Biometrika* **87**, 99–112.
- SCOTT, J. G. & BERGER, J. O. (2006). An exploration of aspects of Bayesian multiple testing. *J. Statist. Plann. Inference* **136**, 2144–62.
- SCOTT, J. G. & BERGER, J. O. (2008). Multiple testing, empirical Bayes, and the variable-selection problem. Discussion Paper 2008-10, Duke University Department of Statistical Science.

SCOTT, J. G. & CARVALHO, C. M. (2008). Feature-inclusion stochastic search for Gaussian graphical models. *J. Comput. Graph. Stat.* To appear.

SUN, D. & BERGER, J. O. (2007). Objective priors for the multivariate normal model. In *Bayesian Statistics VIII*, Eds. J. Bernardo, M. Bayarri, J. Berger, A. Dawid, D. Heckerman, A. Smith & M. West, Oxford University Press, to appear.

ZELLNER, A. (1986). On assessing prior distributions and Bayesian regression analysis with g -prior distributions. In *Bayesian Inference and Decision Techniques: Essays in Honor of Bruno de Finetti*. Elsevier, pp. 233–43.

ZELLNER, A. & SIOW, A. (1980). Posterior odds ratios for selected regression hypotheses. In *Bayesian Statistics: Proceedings of the First International Meeting held in Valencia*, pp. 585–603.

A PROOF OF LEMMA 4.3

Proof. We start by assuming, without loss of generality, that z is a neighbor of all k variables in \mathbf{y} . This can always be achieved by marginalizing over the non-neighbours of z , an operation that changes the conditional independence pattern of \mathbf{y} but maintains the decomposability of the graph (Cowell et al., 1999, Chapter 7) and preserves the hyper-inverse Wishart distribution on the covariance matrix of the remaining vertices (Letac & Massam, 2007). We are now able to apply the results quoted below. Let $\Omega = \Sigma^{-1} = \Phi'\Phi$ be partitioned as:

$$\begin{pmatrix} \Omega_{zz} & \Omega_{zy} \\ \Omega_{yz} & \Omega_{yy} \end{pmatrix} = \begin{pmatrix} \Phi'_{zz} & 0 \\ \Phi'_{yz} & \Phi'_{yy} \end{pmatrix} \begin{pmatrix} \Phi_{zz} & \Phi_{zy} \\ 0 & \Phi_{yy} \end{pmatrix}.$$

Recall that $(\Sigma_{zz} - \Sigma_{zy}\Sigma_{yy}^{-1}\Sigma_{yz})^{-1} = \Sigma_{z|y}^{-1} = \Omega_{zz} = \Phi'_{zz}\Phi_{zz}$ and $\Sigma_{zy}\Sigma_{yy}^{-1} = -\Phi_{zz}^{-1}\Phi_{zy} = \Omega_{zz}^{-1}\Omega_{zy}$. If $(\Sigma|G) \sim \text{HIW}_G(b, D)$, properties of the Cholesky decomposition of the hyper-inverse Wishart as defined in Roverato (2000) and Atay-Kayis & Massam (2005) allow us to write $\Psi = \Phi T^{-1}$ where $D^{-1} = T'T$ with

$$\Psi_{zz}^2 \sim \text{GA}\left(\frac{b+k}{2}, \frac{1}{2}\right) \quad \text{and} \quad \Psi_{zy} \sim \text{N}(0, I). \quad (17)$$

Notice that this is true regardless of the order in which the nodes are listed as z is a neighbour of all variables in \mathbf{y} and therefore all elements in Ψ_{zy} are unconstrained. It is straightforward to see from (17) that

$$\Omega_{zz} = \Phi'_{zz}\Phi_{zz} = \Phi_{zz}^2 = (\Psi_{zz}T_{zz})^2 \sim \text{GA}\left(\frac{b+k}{2}, \frac{T_{zz}^{-2}}{2}\right),$$

so that

$$\Sigma_{z|\mathbf{y}}^{-1} \sim \text{GA} \left(\frac{b+k}{2}, \frac{D_{z|\mathbf{y}}}{2} \right),$$

which proves part (i) of the Lemma. Turning the focus to the form of $\Gamma = \Sigma_{z\mathbf{y}}\Sigma_{\mathbf{y}\mathbf{y}}^{-1} = -\Phi_{zz}^{-1}\Phi_{z\mathbf{y}}$ and writing it as a function of Ψ and T , we get

$$\gamma_i = - \left(\frac{T_{y_i z}}{T_{zz}} + \frac{1}{\Psi_{zz} T_{zz}} \sum_{j=1}^i \Psi_{zy_j} T_{y_j y_i} \right). \quad (18)$$

Given Φ_{zz} , this is just a linear combination of independent standard normals, so that

$$(\Gamma | \Phi_{zz}) \sim \text{N} \left(-T_{zz}^{-1} T_{z\mathbf{y}}, \frac{1}{\Phi_{zz}^2} T'_{\mathbf{y}\mathbf{y}} T_{\mathbf{y}\mathbf{y}} \right) \quad (19)$$

$$(\Sigma_{z\mathbf{y}}\Sigma_{\mathbf{y}\mathbf{y}}^{-1} | \Sigma_{z|\mathbf{y}}^{-1}) \sim \text{N} (D_{z\mathbf{y}} D_{\mathbf{y}\mathbf{y}}^{-1}, \Sigma_{z|\mathbf{y}} D_{\mathbf{y}\mathbf{y}}^{-1}), \quad (20)$$

proving part (ii) of the Lemma. □