

# COPULA GAUSSIAN GRAPHICAL MODELS AND THEIR APPLICATION TO MODELING FUNCTIONAL DISABILITY DATA<sup>1</sup>

BY ADRIAN DOBRA AND ALEX LENKOSKI

*University of Washington and Heidelberg University*

We propose a comprehensive Bayesian approach for graphical model determination in observational studies that can accommodate binary, ordinal or continuous variables simultaneously. Our new models are called copula Gaussian graphical models (CGGMs) and embed graphical model selection inside a semiparametric Gaussian copula. The domain of applicability of our methods is very broad and encompasses many studies from social science and economics. We illustrate the use of the copula Gaussian graphical models in the analysis of a 16-dimensional functional disability contingency table.

**1. Introduction.** The determination of conditional independence relationships through graphical models is a key component of the statistical analysis of observational studies. A pertinent example we will focus on in this paper is a functional disability data set extracted from the “analytic” data file for the National Long Term Care Survey (NLTC) created by the Center of Demographic Studies at Duke University. Each observed variable is binary and corresponds to a measure of disability defined by an activity of daily living. This contingency table cross-classifies information on elderly aged 65 and above pooled across four survey waves, 1982, 1984, 1989 and 1994—see [Manton, Corder and Stallard \(1993\)](#) for more details. The 16 dimensions of this table correspond to six activities of daily living (ADLs) and ten instrumental activities of daily living (IADLs). Specifically, the ADLs relate to hygiene and personal care: eating (ADL1), getting in/out of bed (ADL2), getting around inside (ADL3), dressing (ADL4), bathing (ADL5) and getting to the bathroom or using a toilet (ADL6). The IADLs relate to activities needed to live without dedicated professional care: doing heavy house work (IADL1), doing light house work (IADL2), doing laundry (IADL3), cooking (IADL4), grocery shopping (IADL5), getting about outside (IADL6), travelling (IADL7), managing money (IADL8), taking medicine (IADL9) and telephoning (IADL10). For each ADL/IADL measure, subjects were classified as being either healthy (level 1) or disabled (level 2) on that measure. The methodology we develop in this paper allows us to determine the complex pattern of conditional

---

Received March 2009; revised June 2010.

<sup>1</sup>Supported in part by NIH Grant R01 HL092071.

*Key words and phrases.* Bayesian inference, Gaussian graphical models, latent variable model, Markov chain Monte Carlo.

associations that exist among the 16 daily living activities. This represents a critical issue that was left unexplored in previous analyses of this data set [Erosheva, Fienberg and Joutard (2007); Fienberg et al. (2010)].

In fact, the domain of applicability of our methods is not restricted to contingency tables. Since multivariate data sets arising from social science or economics typically contain variables of many types, our goal is to develop an approach to graphical model determination that is broad enough to be applicable to any study that involves a mixture of binary, ordinal and continuous variables.

Most of the research efforts in the graphical models literature have been focused on multivariate normal models or on log-linear models—see, for example, the monographs of Lauritzen (1996) and Whittaker (1990). These models relate to data sets that contain exclusively continuous or categorical variables. CG distributions [Lauritzen (1996)] constitute the basis of a class of graphical models for mixed variables, but they impose an overly restrictive assumption: the conditional distribution of the continuous variables given the discrete variables must be multivariate normal. As such, the three main classes of graphical models are too restrictive to be widely applicable to social science or economics studies.

Copulas [Nelsen (1999)] provide the theoretical framework in which multivariate associations can be modeled separately from the univariate distributions of the observed variables. Genest and Neslehová (2007) advocate the use of copulas when modeling multivariate distributions involving discrete variables. In this paper we employ the Gaussian copula and further require conditional independence constraints on the inverse of its correlation matrix. The resulting models are called copula Gaussian graphical models (CGGMs) because they only impose a multivariate normal assumption for a set of latent variables which are in a one-to-one correspondence with the set of observed variables. A related approach for inference in Gaussian copulas has been developed by Pitt, Chan and Kohn (2006). Their framework involves parametric models for Gaussian copulas and the univariate marginal distributions of the observed variables. We treat these marginal distributions as nuisance parameters and focus on the determination of graphical models.

The structure of the paper is as follows. In Section 2 we formally introduce Gaussian graphical models (GGMs) and describe a Bayesian framework for inference in this class of models. In Section 3 we discuss modeling aspects related to binary and ordinal variables. In Section 4 we show how to extend GGMs to represent conditional independence associations in a latent variables space. We also present a Bayesian model averaging approach for graph identification and estimation in CGGMs. In Section 5 we analyze the NLTCs functional disability data together with another six-dimensional contingency table using CGGMs. We discuss our proposed methodology in Section 6.

**2. Gaussian graphical models.** We let  $X = X_V$ ,  $V = \{1, 2, \dots, p\}$ , be a random vector with a joint distribution  $p(X_V)$ . The conditional independence

relationships among  $\{X_v : v \in V\}$  under  $p(X_V)$  can be summarized in a graph  $G = (V, E)$ , where each vertex  $v \in V$  corresponds with a random variable  $X_v$  and  $E \subset V \times V$  are undirected edges [Whittaker (1990)]. Here “undirected” means that  $(v_1, v_2) \in E$  is equivalent with  $(v_2, v_1) \in E$ .

The absence of an edge between  $X_{v_1}$  and  $X_{v_2}$  corresponds with the conditional independence of these two random variables given the remaining variables under  $p(X_V)$  and is denoted by

$$(2.1) \quad X_{v_1} \perp\!\!\!\perp X_{v_2} \mid X_{V \setminus \{v_1, v_2\}}.$$

This is called the pairwise Markov property relative to  $G$ , which in turn implies the local as well as the global Markov properties relative to  $G$  [Lauritzen (1996)].

We denote by  $\mathcal{G}_V$  the set of all  $2^{p(p-1)/2}$  undirected graphs with vertices  $V$ . Since  $\mathcal{G}_V$  contains many graphs even for relatively small values of  $p$ , it cannot be enumerated and has to be visited using stochastic search methods [Madigan and York (1995); Jones et al. (2005); Lenkoski and Dobra (2010)]. Such algorithms move through  $\mathcal{G}_V$  using neighborhood sets  $\text{nbr}(G) \subset \mathcal{G}_V$  for  $G \in \mathcal{G}_V$ . The neighborhood of a graph  $G \in \mathcal{G}_V$  is comprised of all the graphs obtained from  $G$  by adding or deleting one edge. These neighborhood sets are symmetric and link any two graphs through a path of graphs such that two consecutive graphs on this path are neighbors of each other. We remark that the neighborhood sets associated with  $\mathcal{G}_V$  contain the same number of graphs  $p(p-1)/2$ .

Furthermore, we assume that  $X = X_V$  follows a  $p$ -dimensional multivariate normal distribution  $N_p(0, K^{-1})$  with precision matrix  $K = (K_{v_1, v_2})_{1 \leq v_1, v_2 \leq p}$ . We let  $x^{(1:n)} = (x^{(1)}, \dots, x^{(n)})^T$  be the observed data of  $n$  independent samples of  $X$ . The likelihood function is proportional to

$$(2.2) \quad p(x^{(1:n)} \mid K) \propto (\det K)^{n/2} \exp\left\{-\frac{1}{2} \langle K, U \rangle\right\},$$

where  $U = \sum_{j=1}^n x^{(j)} x^{(j)T}$ , and  $\langle A, B \rangle = \text{tr}(A^T B)$  denotes the trace inner product. We assume that the data have been centered and scaled, so that the sample mean of each  $X_v$  is zero and its sample variance is one.

A graphical model  $G = (V, E)$  for  $N_p(0, K^{-1})$  is called a Gaussian graphical model (GGM) and is constructed by constraining some of the off-diagonal elements of  $K$  to zero. For example, the pairwise Markov property (2.1) holds if and only if  $K_{v_1, v_2} = 0$ . This implies that the edges of  $G$  correspond with the off-diagonal nonzero elements of  $K$ , that is,  $E = \{(v_1, v_2) \mid K_{v_1, v_2} \neq 0, v_1 \neq v_2\}$ . Given  $G$ , the precision matrix  $K$  is constrained to the cone  $P_G$  of symmetric positive definite matrices with entries  $K_{v_1, v_2}$  equal to zero for all  $(v_1, v_2) \notin E, v_1 \neq v_2$ .

We consider a  $G$ -Wishart prior  $W_G(\delta, D)$  for  $K$  with density

$$(2.3) \quad p(K \mid G) = \frac{1}{I_G(\delta, D)} (\det K)^{(\delta-2)/2} \exp\left\{-\frac{1}{2} \langle K, D \rangle\right\},$$

with respect to the Lebesgue measure on  $P_G$  [Roverato (2002); Atay-Kayis and Massam (2005); Letac and Massam (2007)]. The normalizing constant  $I_G(\delta, D)$

is finite provided  $\delta > 2$  and  $D$  is positive definite [Diaconis and Ylvisaker (1979)]. If  $G$  is the complete graph with  $p$  vertices (i.e., there are no missing edges),  $W_G(\delta, D)$  reduces to the Wishart distribution  $W_p(\delta, D)$ , hence, its normalizing constant is

$$(2.4) \quad I_G(\delta, D) = 2^{(\delta+p-1)p/2} \Gamma_p\{(\delta + p - 1)/2\} (\det D)^{-(\delta+p-1)/2},$$

where  $\Gamma_p(a) = \pi^{p(p-1)/4} \prod_{i=0}^{p-1} \Gamma(a - \frac{i}{2})$  for  $a > (p - 1)/2$  [Muirhead (2005)]. If  $G$  is decomposable,  $I_G(\delta, D)$  is explicitly calculated [Roverato (2002)]. For nondecomposable graphs, the Monte Carlo method of Atay-Kayis and Massam (2005) can be used to numerically approximate  $I_G(\delta, D)$  in a fast and accurate manner.

Throughout this paper we set the prior parameters for  $K$  to  $\delta = 3$  and  $D = I_p$ , the  $p$ -dimensional identity matrix. From equations (2.2) and (2.3) we see that the interpretation of this prior is that the components of  $X$  are independent apriori and that the “weight” of the prior is equivalent to one observed sample.

The  $G$ -Wishart prior is conjugate to the likelihood (2.2), thus, the posterior distribution of  $K$  given  $G$  is  $W_G(\delta + n, D + U)$ , that is,

$$p(K|x^{(1:n)}, G) = \frac{1}{I_G(\delta + n, D + U)} (\det K)^{(\delta+n-2)/2} \exp\left\{-\frac{1}{2}\langle K, D + U \rangle\right\}.$$

Given  $K \in P_G$ , the regression of  $X_v$  on the remaining elements of  $X$  depends only on the neighbors of  $v$  in  $G$ :

$$(2.5) \quad p(X_v | X_{V \setminus \{v\}} = x_{V \setminus \{v\}}, K) = N\left(-\sum_{v' \in bd_G(v)} \frac{K_{v,v'}}{K_{v,v}} x_{v'}, \frac{1}{K_{v,v}}\right),$$

where  $bd_G(v) = \{v' \in V : (v, v') \in E\}$ .

The Cholesky decomposition of a matrix  $K \in P_G$  is  $K = \phi^T \phi$ , where  $\phi$  is an upper triangular matrix with  $\phi_{v,v} > 0, v \in V$ . Roverato (2002) proved that the set  $\nu(G)$  of the free elements of  $\phi$  consists of the diagonal elements together with the elements that correspond with the edges of  $G$ , that is,

$$\nu(G) = \{(v_1, v_1) : v_1 \in V\} \cup \{(v_1, v_2) : v_1 < v_2 \text{ and } (v_1, v_2) \in E\}.$$

Once the free elements of  $\phi$  are known, the remaining elements are also known. More specifically, we have  $\phi_{1,v_2} = 0$  if  $v_2 \geq 2$  and  $(1, v_2) \notin E$ . We also have

$$\phi_{v_1, v_2} = -\frac{1}{\phi_{v_1, v_1}} \sum_{v=1}^{v_1-1} \phi_{v, v_1} \phi_{v, v_2}$$

for  $2 \leq v_1 < v_2$  and  $(v_1, v_2) \notin E$ . The determination of the elements of  $\phi$  that are not free based on the elements of  $\phi$  that are free is called the completion of  $\phi$  with respect to  $G$  [Roverato (2002); Atay-Kayis and Massam (2005)]. It is useful to remark that the free elements of  $\phi$  fully determine the matrix  $K$ . The development

of our framework involves the Jacobian of the transformation that maps  $K \in P_G$  to the free elements of  $\phi$  [Roverato (2002)]:

$$J(K \rightarrow \phi) = 2^p \prod_{v=1}^p \phi_{v,v}^{d_v^G+1},$$

where  $d_v^G$  is the number of elements in  $bd_G(v) \cap \{v + 1, \dots, p\}$ .

**3. Incorporating binary and ordinal categorical variables.** A variable  $X_v$  that takes a finite number of ordinal values  $\{1, 2, \dots, d_v\}$ , with  $d_v \geq 2$ , is incorporated in our modeling framework by introducing a continuous latent variable  $Z_v$  underlying  $X_v$ —see, for example, Muthén (1984). We denote by  $\{x_v^{(1)}, \dots, x_v^{(n)}\}$  the observed samples associated with  $X_v$ . The samples from  $Z_v$  are denoted by  $\{z_v^{(1)}, \dots, z_v^{(n)}\}$ . Typically the relationship between  $X_v$  and its surrogate  $Z_v$  is expressed through some thresholds  $\tau_v = (\tau_{v,0}, \tau_{v,1}, \dots, \tau_{v,w_v})$  with  $-\infty = \tau_{v,0} < \tau_{v,1} < \dots < \tau_{v,w_v} = \infty$ . Formally, we set [Dunson (2006)]

$$(3.1) \quad x_v^{(j)} = \sum_{l=1}^{w_v} l \times \mathbf{1}_{\{\tau_{v,l-1} < z_v^{(j)} \leq \tau_{v,l}\}}, \quad j = 1, 2, \dots, n.$$

This model is identifiable if the value of  $\tau_{v,1}$  is fixed at a certain value. We follow an idea originally suggested by Hoff (2007) that does not explicitly involve the thresholds  $\tau_v$ . This approach is based on the remark that the relationship between the observed and latent samples satisfies the constraints

$$(3.2) \quad x_v^{(j_1)} < x_v^{(j_2)} \Rightarrow z_v^{(j_1)} < z_v^{(j_2)}, \quad z_v^{(j_1)} < z_v^{(j_2)} \Rightarrow x_v^{(j_1)} \leq x_v^{(j_2)}$$

for  $1 \leq j_1 \neq j_2 \leq n$ . We see that if  $X_v$  and  $Z_v$  are related as in (3.1), then (3.2) holds. If (3.2) holds, then (3.1) also holds by choosing  $\tau_{v,l} = \max\{z_v^{(j)} : x_v^{(j)} = l\}$  for  $l = 1, \dots, w_v - 1$ . It follows that, given the observed data  $x^{(1:n)}$ , the latent samples  $z^{(1:n)} = (z^{(1)}, z^{(2)}, \dots, z^{(n)})$  are constrained to belong to the set

$$A(x^{(1:n)}) = \{z^{(1:n)} \in \mathbf{R}^{n \times p} : L_v^j(z^{(1:n)}) < z_v^{(j)} < U_v^j(z^{(1:n)})\},$$

where

$$(3.3) \quad \begin{aligned} L_v^j(z^{(1:n)}) &= \max\{z_v^{(k)} : x_v^{(k)} < x_v^{(j)}\}, \\ U_v^j(z^{(1:n)}) &= \min\{z_v^{(k)} : x_v^{(j)} < x_v^{(k)}\}. \end{aligned}$$

If the value  $x_v^{(j)}$  is missing from the observed data, we define  $L_v^j(z^{(1:n)}) = -\infty$  and  $U_v^j(z^{(1:n)}) = \infty$ .

**4. Copula Gaussian graphical models.** We assume that an observed variable  $X_v$  can be binary, categorical with ordered categories, count or continuous. We denote by  $F_v$  the univariate distribution of  $X_v$  and by  $F_v^{-1}$  the pseudo-inverse of  $F_v$ . Given a precision matrix  $K$ , we model the joint distribution of  $X = X_V$  as follows [see also Hoff (2007)]:

$$\begin{aligned}
 Z_V &\sim N_p(0, K^{-1}), \\
 \tilde{Z}_v &= Z_v / (K^{-1})_{v,v}^{1/2}, \quad v \in V, \\
 X_v &= F_v^{-1}(\Phi(\tilde{Z}_v)), \quad v \in V.
 \end{aligned}
 \tag{4.1}$$

In (4.1) the joint distribution of the latent variables is multivariate normal  $\tilde{Z} = \tilde{Z}_V \sim N_p(0, \Upsilon(K))$ , where  $\Upsilon(K)$  is a correlation matrix with entries

$$\Upsilon_{v_1, v_2}(K) = \frac{(K^{-1})_{v_1, v_2}}{\sqrt{(K^{-1})_{v_1, v_1} (K^{-1})_{v_2, v_2}}}.
 \tag{4.2}$$

The joint distribution  $F$  of  $X = X_V$  is subsequently a function of the correlation matrix  $\Upsilon(K)$  and the univariate distributions  $F_v$  of  $X_v$ :

$$\begin{aligned}
 p(X_1 \leq x_1, \dots, X_p \leq x_p) &= F(x_1, \dots, x_p | \Upsilon(K), F_1, \dots, F_p), \\
 &= C(F_1(x_1), \dots, F_p(x_p) | \Upsilon(K)),
 \end{aligned}$$

where

$$C(u_1, \dots, u_p | \Upsilon') = \Phi_p(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_p) | \Upsilon') : [0, 1]^p \rightarrow [0, 1]
 \tag{4.3}$$

is the Gaussian copula with  $p \times p$  correlation matrix  $\Upsilon'$  [Nelsen (1999)]. Here  $\Phi(\cdot)$  represents the CDF of the standard normal distribution and  $\Phi_p(\cdot | \Upsilon)$  is the CDF of  $N_p(0, \Upsilon)$ .

We avoid the need to formally make assumptions regarding the parametric representation of  $\{F_v : v \in V\}$ , which could be a daunting task for most real world data sets, by treating their marginal distributions as nuisance parameters. Moreover, we reduce our model parameters to the correlation matrix of the Gaussian copula (4.3). This means that we focus on the joint distribution of the latent variables  $\tilde{Z}_V$  whose relationships with the observed variables  $X_V$  are given by (4.1). Since  $F_v^{-1}(\cdot)$  and  $\Phi(\cdot)$  are nondecreasing, (4.1) implies (3.2) which does not depend on the marginal distributions  $\{F_v : v \in V\}$ . The converse is also true: if the relationship (3.2) between the observed and latent samples holds, then (4.1) also holds by replacing  $F_v$  with the empirical distribution of  $X_v$ .

As suggested by Hoff (2007), inference in the latent variables space can be performed by substituting the observed data  $x^{(1:n)}$  with the event  $\mathcal{D} = \{z^{(1:n)} \in A(x^{(1:n)})\}$ . We write the likelihood function as

$$p(x^{(1:n)} | K, \{F_v : v \in V\}) = p(\mathcal{D} | K) p(x^{(1:n)} | \mathcal{D}, \Upsilon(K), \{F_v : v \in V\}).$$

In this decomposition  $p(\mathcal{D}|K)$  is the only part of the observed data likelihood that is relevant for making inference on  $K$ . Furthermore,  $p(\mathcal{D}|K)$  does not depend on  $\{F_v : v \in V\}$ . Hoff (2007) calls  $p(\mathcal{D}|K)$  the extended rank likelihood and constructs a Gibbs sampler with stationary distribution

$$(4.4) \quad p(K|\mathcal{D}) \propto p(\mathcal{D}|K)p(K),$$

where  $K$  follows a Wishart prior distribution  $W_p(\delta, D)$ .

We are interested in modeling the conditional independence relationships among the latent variables  $Z = Z_V$  using Gaussian graphical models. We go one step further compared to Hoff (2007) and impose zero constraints in the precision matrix  $K$  according to a graph  $G$ . We refer to the graphical models constructed in the latent space as copula Gaussian graphical models (CGGMs). The inference approach described in Hoff (2007) is equivalent to reducing the set of candidate graphs to only one graph. This graph is the full graph in which all the edges are present and none of the off-diagonal elements of  $K$  are constrained to zero.

The Markov properties associated with a CGGM are guaranteed to translate into Markov properties for the observed variables if all the marginals  $\{F_v : v \in V\}$  are continuous [Liu, Lafferty and Wasserman (2009)]. The presence of some discrete observed variables might induce additional dependencies among the  $X$ 's that are not modeled in a CGGM, but such dependencies can be regarded as having a secondary relevance since they emerge from the marginals  $\{F_v : v \in V\}$ . The conditional independence graphs for the latent variables could contain edges then that do not necessarily correspond with conditional independence relationships in the observed variables space. Conversely, there might exist conditional independence relationships among the observed variables that are not represented in conditional independence graphs that involve latent variables.

4.1. *Bayesian inference in copula Gaussian graphical models.* Let  $G \in \mathcal{G}_V$  be a graph defining a CGGM. The joint posterior distribution of  $K \in P_G$  and the graph  $G$  is given by

$$(4.5) \quad p(K, G|\mathcal{D}) \propto p(\mathcal{D}|K)p(K|G)p(G).$$

The prior distribution of  $K$  conditional on  $G$  is  $G$ -Wishart  $W_G(\delta, D)$  and the prior distribution over  $\mathcal{G}_V$  is uniform, that is,  $p(G) \propto 1$ . Other choices of priors on the graphs space  $\mathcal{G}_V$  take into consideration the implied distribution on the number of edges [Wong, Carter and Kohn (2003)], encourage sparsity [Jones et al. (2005)] or have multiple testing correction properties [Scott and Berger (2006)].

We describe a Markov chain Monte Carlo sampler for the joint distribution (4.5). We consider two strictly positive precision parameters  $\sigma_p$  and  $\sigma_g$  that remain fixed throughout at some small values, for example,  $\sigma_p = \sigma_g = 0.1$ . Given the current state of the chain  $(K^s, G^s)$ , its next state  $(K^{s+1}, G^{s+1})$  is generated by sequentially performing the following updates.

*Step 1: Resample the latent data.* For each  $v \in V$  and  $j \in \{1, 2, \dots, n\}$ , we update the latent value  $z_v^{(j)}$  by sampling from its full conditional distribution. The distribution of  $Z_v$  conditional on  $Z_{V \setminus \{v\}} = z_{V \setminus \{v\}}^{(j)}$  is  $N(\mu_v, \sigma_v^2)$  truncated to the interval  $[L_v^j, U_v^j]$ , where  $\mu_v = -\sum_{v' \in \text{bd}_G(v)} \frac{K_{v,v'}^s}{K_{v,v}^s} z_{v'}^{(j)}$  and  $\sigma_v^2 = \frac{1}{K_{v,v}^s}$ —see (2.5). The bounds  $L_v^j$  and  $U_v^j$  are given in (3.3). The new value of  $z_v^{(j)}$  is obtained by sampling from this truncated normal distribution.

*Step 2: Resample the precision matrix.* We sequentially perturb the free elements  $\{\phi_{v_1, v_2}^s : (v_1, v_2) \in \nu(G^s)\}$  in the Cholesky decomposition  $K^s = (\phi^s)^T \phi^s$  around their current value. Here  $\phi^s$  is upper triangular. We perform a Metropolis–Hastings update of  $K^s$  associated with a diagonal element  $\phi_{v_1, v_1}^s > 0$  by sampling a value  $\gamma$  from a  $N(\phi_{v_1, v_1}^s, \sigma_p^2)$  distribution truncated below at 0, that is,

$$\gamma \sim q(u | \phi_{v_1, v_1}^s) \propto \frac{1}{\sigma_p \Phi(\phi_{v_1, v_1}^s / \sigma_p)} \exp\left(-\frac{(u - \phi_{v_1, v_1}^s)^2}{2\sigma_p^2}\right).$$

We take  $K' = (\phi')^T \phi'$ , where  $\phi'$  is such that its free elements coincide with the free elements of  $\phi^s$ , with the exception of the  $(v_1, v_1)$  element which is set to  $\gamma$ . The elements of  $\phi'$  that are not free are obtained by the completion operation described in Section 2. The acceptance probability of the update of  $K^s$  to  $K'$  is  $\min\{R_p, 1\}$ , where

$$\begin{aligned} R_p &= \frac{p(K' | z^{(1:n)}, G^s)}{p(K^s | z^{(1:n)}, G^s)} \frac{J(K' \rightarrow \phi')}{J(K^s \rightarrow \phi^s)} \frac{q(\phi_{v_1, v_1}^s | \gamma)}{q(\gamma | \phi_{v_1, v_1}^s)}, \\ &= \frac{\Phi(\phi_{v_1, v_1}^s / \sigma_p)}{\Phi(\gamma / \sigma_p)} \left(\frac{\gamma}{\phi_{v_1, v_1}^s}\right)^{\delta + n + d_{v_1}^{G^s} - 1} R'_p. \end{aligned}$$

Here we denote

$$R'_p = \exp\left\{-\frac{1}{2} \left\langle K' - K^s, D + \sum_{j=1}^n z^{(j)} z^{(j)T} \right\rangle\right\}.$$

Next we consider a free off-diagonal element  $\phi_{v_1, v_2}^s$ , where  $v_1 < v_2$  and  $(v_1, v_2) \in \nu(G^s)$ . We sample a candidate value  $\gamma'$  from a  $N(\phi_{v_1, v_2}^s, \sigma_p^2)$  distribution. As before, we take  $K' = (\phi')^T \phi'$ , where  $\phi'$  and  $\phi^s$  have the same free elements with the exception of the  $(v_1, v_2)$  element that has  $\phi'_{v_1 v_2} = \gamma'$ . The remaining nonfree elements of  $\phi'$  are obtained through completion. Due to the symmetry of the proposal distribution and the fact that  $\det K^s = \prod_{v=1}^p (\phi_{v,v}^s)^2 = \prod_{v=1}^p (\phi'_{v,v})^2 = \det K'$ , the candidate matrix  $K'$  is accepted with probability  $\min\{R'_p, 1\}$ .

Since  $K^s \in P_{G^s}$ , the candidate matrix  $K'$  associated with each free element in  $\nu(G^s)$  must also belong to  $P_{G^s}$ . The precision matrix that is obtained after performing all the Metropolis–Hastings updates is  $K^{s+1/2} \in P_{G^s}$ .



*Step 3: Resample the graph.* We consider the Cholesky decomposition  $K^{s+1/2} = (\phi^{s+1/2})^T \phi^{s+1/2}$  where  $\phi^{s+1/2}$  is upper triangular. We randomly choose a pair  $(v_1, v_2)$ ,  $v_1 < v_2$ . If there is no edge between  $v_1$  and  $v_2$  in  $G^s$ , that is,  $(v_1, v_2) \notin v(G^s)$ , we add this edge to  $G^s$  to obtain a candidate graph  $G'$ . This implies  $bd_{G'}(v_1) = bd_{G^s}(v_1) \cup \{v_2\}$ , hence,  $d_{v_1}^{G'} = d_{v_1}^{G^s} + 1$ . Moreover,  $v(G') = v(G^s) \cup \{(v_1, v_2)\}$ . We define an upper diagonal matrix  $\phi'$  such that  $\phi'_{v'_1, v'_2} = \phi^{s+1/2}_{v'_1, v'_2}$  for all  $(v'_1, v'_2) \in v(G^s)$ . The value of  $\phi'_{v_1, v_2}$  is set by sampling from a  $N(\phi^{s+1/2}_{v_1, v_2}, \sigma_g^2)$  distribution. The remaining elements of  $\phi'$  are determined through completion with respect to the graph  $G'$ . We see that  $\phi'$  has one additional free element with respect to  $\phi^{s+1/2}$  whose value was randomly chosen by perturbing the nonfree  $(v_1, v_2)$  element of  $\phi^{s+1/2}$ .

We take the candidate precision matrix  $K' = (\phi')^T \phi' \in P_{G'}$ . Since the dimensionality of the parameter space increases by one, we must make use of the reversible jump Markov chains methodology proposed by Green (1995). We accept the update of  $(K^{s+1/2}, G^s)$  to  $(K', G')$  with probability  $\min\{R_g, 1\}$ , where  $R_g$  is given by

$$\frac{p(z^{(1:n)}|K')p(K'|G')}{p(z^{(1:n)}|K^{s+1/2})p(K^{s+1/2}|G^s)} \frac{|\text{nbnd}(G^s)|}{|\text{nbnd}(G')|} \times \frac{J(K' \rightarrow \phi')}{J(K^{s+1/2} \rightarrow \phi^{s+1/2})} \frac{J(\phi^{s+1/2} \rightarrow \phi')}{(1/(\sigma_g \sqrt{2\pi})) \exp(-(\phi'_{v_1, v_2} - \phi^{s+1/2}_{v_1, v_2})^2 / (2\sigma_g^2))}.$$

We denote by  $|B|$  the number of elements of a set  $B$ . All the graphs in  $\mathcal{G}_V$  have the same number of neighbors, hence,  $|\text{nbnd}(G^s)| = |\text{nbnd}(G')| = p(p - 1)/2$ . Since the free elements of  $\phi'$  are the free elements of  $\phi^{s+1/2}$  and  $\phi'_{v_1, v_2}$ , the Jacobian of the transformation from  $\phi^{s+1/2}$  to  $\phi'$  is equal to 1, that is,  $J(\phi^{s+1/2} \rightarrow \phi') = 1$ . Moreover,  $\phi^{s+1/2}$  and  $\phi'$  have the same elements on the main diagonal and are upper triangular, therefore,  $\det K^{s+1/2} = \det K'$ . We also have

$$\frac{J(K' \rightarrow \phi')}{J(K^{s+1/2} \rightarrow \phi^{s+1/2})} = \frac{(\phi'_{v_1, v_1})^{d_{v_1}^{G'} + 1}}{(\phi^{s+1/2}_{v_1, v_1})^{d_{v_1}^{G^s} + 1}} = \phi^{s+1/2}_{v_1, v_1}.$$

It follows that  $R_g$  is equal to

$$\sigma_g \sqrt{2\pi} \phi^{s+1/2}_{v_1, v_1} \frac{I_{G^s}(\delta, D)}{I_{G'}(\delta, D)} \times \exp \left\{ -\frac{1}{2} \left\langle K' - K^{s+1/2}, D + \sum_{j=1}^n z^{(j)} z^{(j)T} \right\rangle + \frac{(\phi'_{v_1, v_2} - \phi^{s+1/2}_{v_1, v_2})^2}{2\sigma_g^2} \right\}.$$

Now we examine the case when there is an edge between  $v_1$  and  $v_2$  in  $G^s$ . We delete this edge from  $G^s$  to obtain a candidate graph  $G'$ . We have  $bd_{G'}(v_1) =$

$bd_{G^s}(v_1) \setminus \{v_2\}$ , hence,  $d_{v_1}^{G'} = d_{v_1}^{G^s} - 1$  and  $v(G') = v(G^s) \setminus \{(v_1, v_2)\}$ . We define an upper diagonal matrix  $\phi'$  such that  $\phi'_{v'_1, v'_2} = \phi^{s+1/2}_{v'_1, v'_2}$  for all  $(v'_1, v'_2) \in v(G')$ . The  $(v_1, v_2)$  element is free in  $\phi^{s+1/2}$ , but it is no longer free in  $\phi'$ . The nonfree elements of  $\phi'$  are obtained by completion with respect to the graph  $G'$ . As before, we take  $K' = (\phi')^T \phi' \in P_{G'}$ . The dimensionality of the parameter space decreases by 1 as we move from  $\phi^{s+1/2}$  to  $\phi'$ . We obtain that the acceptance probability of the update from  $(K^{s+1/2}, G^s)$  to  $(K', G')$  is  $\min\{R'_g, 1\}$ , where  $R'_g$  is equal to

$$(\sigma_g \sqrt{2\pi} \phi^{s+1/2}_{v_1, v_1})^{-1} \frac{I_{G^s}(\delta, D)}{I_{G'}(\delta, D)} \times \exp \left\{ -\frac{1}{2} \left\langle K' - K^{s+1/2}, D + \sum_{j=1}^n z^{(j)} z^{(j)T} \right\rangle - \frac{(\phi'_{v_1, v_2} - \phi^{s+1/2}_{v_1, v_2})^2}{2\sigma_g^2} \right\}.$$

The updated graph and the corresponding precision matrix that are obtained at the end of this step are  $G^{s+1}$  and  $K^{s+1}$ , respectively.

We note that our strategy for updating the precision matrix and the graph has some similarities with the work of Giudici and Green (1999). However, they focused exclusively on decomposable graphs and perturbed elements of the covariance matrix  $K^{-1}$  that are either on its main diagonal or correspond to an edge in the graph.

*4.2. Estimation and testing in copula Gaussian graphical models.* In high-dimensional data sets with a small number of observed samples it is likely that the highest posterior probability graph receives only a small (almost zero) posterior probability. Furthermore, changing a few edges in this graph could lead to graphs with comparable posterior probabilities. When model uncertainty is high, Bayesian model averaging becomes key because it avoids the need to perform inference by making an explicit choice about which edges are present or absent in the graphs that underlie the CGGMs. This choice is not desirable since a small sample size means lack of sufficient information. As such, averaging over a large number of graphs is preferable even if prediction is not the final goal.

We let  $\{(G^s, K^s, \Upsilon^s) : s = 1, 2, \dots, S\}$  be samples from the joint distribution (4.5), where  $\Upsilon^s$  is the correlation matrix corresponding with  $K^s$ —see (4.2). These samples can be used to produce Monte Carlo estimates of functions involving the latent variables  $Z$  or the observed variables  $X$ . The posterior probability that two latent variables  $Z_{v_1}$  and  $Z_{v_2}$  are not conditionally independent given  $Z_{V \setminus \{v_1, v_2\}}$  is the posterior inclusion probability of the edge  $(v_1, v_2)$  which is estimated as the proportion of graphs  $G^s$  that contain the edge  $(v_1, v_2)$ .

The posterior expectation of the correlation matrix  $\Upsilon$  is estimated by the mean  $\tilde{\Upsilon} = \frac{1}{S} \sum_{s=1}^S \Upsilon^s$ . A zero element of the correlation matrix  $\Upsilon$  implies the independence of  $Z_{v_1}$  and  $Z_{v_2}$ , which in turn implies the independence of  $X_{v_1}$  and  $X_{v_2}$ . We

can conduct a Bayesian test of independence of  $X_{v_1}$  and  $X_{v_2}$  by considering the interval null hypothesis  $H_{0,\Upsilon}^{v_1,v_2} : |\Upsilon_{v_1,v_2}| < \varepsilon$  with the alternative  $H_{1,\Upsilon}^{v_1,v_2} : |\Upsilon_{v_1,v_2}| \geq \varepsilon$ , where  $\varepsilon > 0$ . Given equal apriori probabilities of the null and alternative hypotheses, the Bayes factor

$$B_{\Upsilon}^{v_1,v_2} = p(H_{1,\Upsilon}^{v_1,v_2} | x^{(1:n)}) / p(H_{0,\Upsilon}^{v_1,v_2} | x^{(1:n)})$$

is estimated as the number of  $\Upsilon_{v_1,v_2}^s$  whose absolute value is above  $\varepsilon$  divided by the number of  $\Upsilon_{v_1,v_2}^s$  whose absolute value is below  $\varepsilon$ .

The CDF of  $X = X_V$  is estimated as

$$\frac{1}{S} \sum_{s=1}^S C(\widehat{F}_1(x_1), \dots, \widehat{F}_p(x_p) | \Upsilon^s),$$

where  $\widehat{F}_v$  is the empirical univariate distribution of  $X_v$ . If each observed variable is discrete and takes values  $\{0, 1, 2, \dots\}$ , their joint probability given  $\Upsilon$  is [Song (2000)]

$$(4.6) \quad p(X_V = x_V | \Upsilon) = \sum_{j_1=0}^1 \dots \sum_{j_p=0}^1 (-1)^{j_1+\dots+j_p} C(u_1^{j_1}(x_1), \dots, u_p^{j_p}(x_p) | \Upsilon),$$

where  $u_v^0(x_v) = \widehat{F}_v(x_v)$  and  $u_v^1(x_v) = \widehat{F}_v(x_v - 1)$ . We define  $u_v^1(0) = 0$ . For example, if  $X_v \in \{0, 1\}$  is a binary random variable, we have  $u_v^0(1) = 1$  and  $u_v^0(0) = u_v^1(1) = \frac{1}{n} \sum_{i=1}^n \delta_{\{x_v^{(i)}=0\}}$ . Here  $\delta_B$  is 1 if  $B$  is true and is 0 otherwise. Thus, the posterior expectation of the joint probability of  $X_V$  is estimated as

$$\tilde{p}(X_V = v_V) = \frac{1}{S} \sum_{s=1}^S p(X_V = x_V | \Upsilon^s).$$

Cramér’s V [Cramér (1946)] is a measure of association between two categorical variables  $X_{v_1}$  and  $X_{v_2}$  that take values in the finite sets  $\mathcal{I}_{v_1}$  and  $\mathcal{I}_{v_2}$ , respectively,

$$(4.7) \quad \rho_{v_1,v_2} = \frac{1}{\min\{|\mathcal{I}_{v_1}|, |\mathcal{I}_{v_2}|\} - 1} \times \sum_{x_{v_1} \in \mathcal{I}_{v_1}} \sum_{x_{v_2} \in \mathcal{I}_{v_2}} \frac{p^2(X_{v_1} = x_{v_1}, X_{v_2} = x_{v_2})}{p(X_{v_1} = x_{v_1})p(X_{v_2} = x_{v_2})} - 1.$$

Cramér’s V always takes values between 0 and 1, but we have  $\rho_{v_1,v_2} = 0$  if and only if  $X_{v_1}$  and  $X_{v_2}$  are independent. The posterior expectation of  $\rho_{v_1,v_2}$  is estimated by calculating the marginal cell value  $p(X_{v_1} = x_{v_1}, X_{v_2} = x_{v_2} | \Upsilon^s)$  of  $p(X_V = x_V | \Upsilon^s)$  for  $s = 1, 2, \dots, S$ , calculating  $\rho_{v_1,v_2}^s$  from (4.7) with respect to  $p(X_{v_1} = x_{v_1}, X_{v_2} = x_{v_2} | \Upsilon^s)$  for  $s = 1, 2, \dots, S$ , then taking the average  $\tilde{\rho}_{v_1,v_2} = \frac{1}{S} \sum_{s=1}^S \rho_{v_1,v_2}^s$ .



women's economic activity. The eight variables are as follows:  $a$ , wife economically active (no, yes);  $b$ , age of wife  $> 38$  (no, yes);  $c$ , husband unemployed (no, yes);  $d$ , child  $\leq 4$  (no, yes);  $e$ , wife's education, high-school+ (no, yes);  $f$ , husband's education, high-school+ (no, yes);  $g$ , Asian origin (no, yes);  $h$ , other household member working (no, yes). The resulting  $2^8$  cross-classification has 165 counts of zero, while 217 cells contain small positive counts smaller than 3. There are quite a few counts larger than 30 or even 50.

Since the sample size is only 665, this table is sparse. Whittaker (1990) argues that higher-order interactions involving more than two variables should not be included in any log-linear model that is fit to this data set. He subsequently studies two log-linear models: the all two-way interaction model whose minimal sufficient statistics are all the 28 two-way marginals and the model whose minimal sufficient statistics are the two-way marginals corresponding with the pairs of variables

$$(5.1) \quad \{fg, ef, dh, dg, cg, cf, ce, bh, be, bd, ag, ae, ad, ac\}.$$

We ran the Markov chain Monte Carlo sampler from Section 4.1 for 250,000 iterations from 100 random starting graphs. The burn-in time was 25,000 iterations. Convergence to the stationary distribution (4.5) is illustrated in Figure 1 that gives the posterior expected number of edges in the CGGM graphs across iterations for each chain. The sampled graphs have on average 16.5 edges which rep-

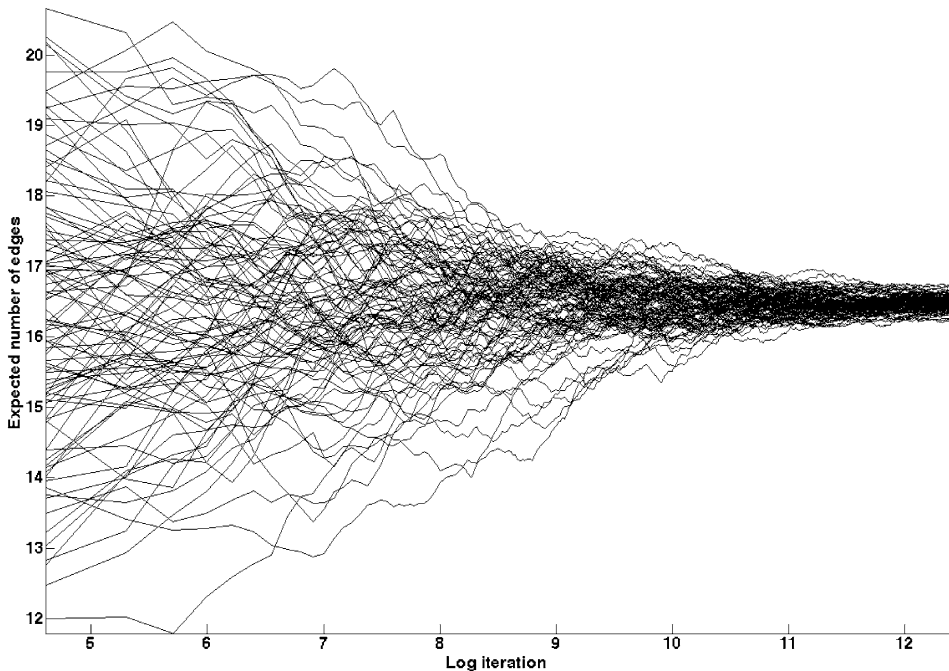


FIG. 1. Estimates of the posterior expected number of edges in the CGGMs for the Rochdale data.

TABLE 2

*Expected cell counts for the top 20 largest counts cells associated with the all two-way interaction log-linear model, Whittaker's log-linear model (5.1), the Copula-Full model and the CGGMs in the Rochdale data. Here 1 stands for no and 2 stands for yes*

Cell	Observed	All two-way	Whittaker	Copula-Full	CGGMs
2 1 1 1 2 2 1 1	57	56.78	52.08	39.43	56.80
2 2 1 1 2 2 1 1	43	44.61	40.97	36.58	47.55
2 2 1 1 1 1 1 1	41	36.40	36.32	30.48	36.12
2 2 1 1 1 2 1 1	37	38.77	36.92	35.33	36.61
2 1 1 2 2 2 1 1	29	33.29	39.06	17.85	32.40
1 1 1 2 2 2 1 1	26	20.36	9.63	9.53	18.03
2 2 1 1 1 2 1 2	26	23.68	22.89	15.67	24.54
2 2 1 1 1 1 1 2	25	28.12	22.52	15.11	27.63
2 1 1 1 1 2 1 1	23	22.73	20.06	26.51	22.76
2 1 1 1 2 1 1 1	22	19.22	16.54	17.15	16.75
2 2 1 1 2 2 1 2	22	22.85	25.41	13.96	24.63
2 1 1 1 1 1 1 1	18	21.54	19.74	21.02	20.85
1 2 1 1 1 1 1 1	17	15.06	16.02	15.13	15.71
1 2 1 1 1 2 1 1	16	14.65	16.28	14.3	12.18
2 2 1 1 2 1 1 1	15	14.96	13.01	17.36	15.07
1 1 1 2 1 2 1 1	13	12.06	6.63	8.46	10.92
2 1 1 2 2 1 1 1	11	7.70	12.40	7.36	8.52
2 1 1 2 1 2 1 1	11	10.50	15.05	11	10.48
1 1 1 2 1 1 2 1	11	8.08	6.72	1.53	6.31

resent approximately 59% of the total number of possible edges. By comparison, the log-linear model (5.1) has 14 minimal sufficient statistics.

In order to show the importance of modeling the conditional independence relationships among the latent variables using graphs, we have also employed the copula estimation approach proposed by Hoff (2007)—see equation (4.4). Hoff’s method is equivalent to starting the Markov chain from Section 4.1 at the full graph and never updating this graph by skipping step 3 of the algorithm. Moreover, updating the precision matrix from step 2 is performed by direct sampling from the Wishart posterior  $W_p(\delta + n, D + \sum_{j=1}^n z^{(j)} z^{(j)T})$ . This simplified Markov chain was run for 25 million iterations and henceforth is called the Copula-Full model.

We compare the expected cell counts of the all two-way interaction log-linear model, the log-linear model (5.1), the Copula-Full model and the CGGMs. Table 2 shows the cells containing the 20 largest observed counts together with their corresponding estimates. It is remarkable that the CGGMs perform as well as the all two-way interaction model for the largest cell count 57. The squared errors between the observed counts and the expected cell counts for all the 256 cells in the table are the following: 284.79 for the all two-way interaction model, 407.04 for the CGGMs, 905.78 for the model (5.1) and 1919.15 for the Copula-Full model.

TABLE 3

*Estimated correlations (elements under the main diagonal) and posterior inclusion probabilities of edges (elements above the main diagonal) associated with the CGGMs in the Rochdale data*

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	—	0.93	0.67	0.92	0.32	0.42	1	0.26
<i>b</i>	0.15	—	0.27	1	0.88	0.29	0.70	0.96
<i>c</i>	-0.52	-0.02	—	0.29	0.91	0.35	0.85	0.25
<i>d</i>	-0.46	-0.79	0.19	—	0.37	0.59	0.66	0.50
<i>e</i>	0.30	-0.28	-0.48	0.12	—	0.98	0.58	0.17
<i>f</i>	0.22	-0.11	-0.35	0.04	0.46	—	0.82	0.22
<i>g</i>	-0.71	-0.31	0.57	0.51	-0.34	-0.37	—	0.32
<i>h</i>	0.12	0.63	0.01	-0.54	-0.19	-0.10	-0.18	—

In Table 3 we show the pairwise correlations  $\Upsilon_{v_1, v_2}$  and the posterior inclusion probabilities of edges  $(v_1, v_2)$  for any two latent variables  $Z_{v_1}$  and  $Z_{v_2}$  as estimated using the CGGMs. In Table 4 we give the estimates of the pairwise correlations  $\Upsilon_{v_1, v_2}$  obtained using the Copula-Full model. We see that the absolute values of these estimates are significantly smaller than corresponding absolute values of the CGGMs estimates. We show the dependence structure of the observed variables in Tables 5 and 6. We give the posterior means of Cramér’s  $V$   $\rho_{v_1, v_2}$  and estimates of the posterior probabilities  $p(H_{1, \rho}^{v_1, v_2} | x^{(1:n)})$  with  $H_{1, \rho}^{v_1, v_2} : \rho_{v_1, v_2} > 0.1$ . By contrasting the estimates obtained using CGGMs and the Copula-Full model, we clearly see that conditioning on the full graph is quite disadvantageous: the Cramér’s  $V$  associations are severely underestimated and, subsequently, all the posterior probabilities  $p(H_{1, \rho}^{v_1, v_2} | x^{(1:n)})$  are almost zero under the full graph. The CGGMs take every possible graph into account and the corresponding estimates are produced by Bayesian model averaging across all graphs. This leads to more appropriate results as evidenced in Tables 3–6.

TABLE 4

*Estimated correlations (elements under the main diagonal) associated with the Copula-Full model in the Rochdale data*

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	—							
<i>b</i>	0.08	—						
<i>c</i>	-0.17	-0.02	—					
<i>d</i>	-0.20	-0.35	0.06	—				
<i>e</i>	0.15	-0.15	-0.15	0.05	—			
<i>f</i>	0.10	-0.06	-0.13	0.02	0.24	—		
<i>g</i>	-0.18	-0.08	0.13	0.13	-0.09	-0.11	—	
<i>h</i>	0.05	0.27	0.01	-0.18	-0.08	-0.06	-0.04	—

TABLE 5

*Estimated Cramér's V associations (elements under the main diagonal) and posterior probabilities  $p(H_{1,\rho}|x^{(1:n)})$  (elements above the main diagonal) associated with the CGGMs in the Rochdale data*

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	—	0	0.19	0.22	0	0	0.83	0
<i>b</i>	0.01	—	0	1	0	0	0	0.94
<i>c</i>	0.08	0	—	0	0	0	0.42	0
<i>d</i>	0.08	0.24	0.01	—	0	0	0.07	0
<i>e</i>	0.04	0.03	0.05	0.01	—	0.35	0	0
<i>f</i>	0.02	0.01	0.03	0	0.09	—	0	0
<i>g</i>	0.12	0.02	0.09	0.07	0.02	0.03	—	0
<i>h</i>	0.01	0.14	0	0.06	0.01	0	0	—

Whittaker (1990), page 282, argues that the strongest pairwise interaction in the Rochdale data is  $(b, d)$ , followed by  $(b, h)$ ,  $(e, f)$  and  $(a, g)$ . In Table 3 we see that the top four posterior inclusion probabilities in the CGGMs are as follows: 1 for  $(b, d)$ , 0.96 for  $(b, h)$ , 0.98 for  $(e, f)$  and 1 for  $(a, g)$ . The strongest associations in the observed variables space as measured by Cramér's V are the following:  $(b, d)$ ,  $(b, h)$ ,  $(a, g)$ ,  $(e, f)$  and  $(c, g)$ . The interaction between  $c$  and  $g$  is also present in the log-linear model (5.1).

Of particular interest is the determination of the factors that influence variable  $a$ —the wife's economic activity. From Table 5 we see that variables  $c$ ,  $d$  and  $g$  are the only variables with a strictly positive posterior probability that their Cramér's V association with variable  $a$  is greater than 0.1. The largest Cramér's V association is  $\tilde{\rho}_{a,g} = 0.12$ , followed by  $\tilde{\rho}_{a,c} = 0.08$  and  $\tilde{\rho}_{a,d} = 0.08$ . The corresponding estimated correlations from Table 3 show a negative relationship between  $a$  and

TABLE 6

*Estimated Cramér's V associations (elements under the main diagonal) and posterior probabilities  $p(H_{1,\rho}|x^{(1:n)})$  (elements above the main diagonal) associated with the Copula-Full model in the Rochdale data*

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>
<i>a</i>	—	0	0	0	0	0	0	0
<i>b</i>	0	—	0	0	0	0	0	0
<i>c</i>	0.01	0	—	0	0	0	0	0
<i>d</i>	0.02	0.04	0	—	0	0	0	0
<i>e</i>	0.01	0.01	0.01	0	—	0	0	0
<i>f</i>	0.01	0	0.01	0	0.03	—	0	0
<i>g</i>	0.01	0	0	0	0	0	—	0
<i>h</i>	0	0.03	0	0	0	0	0	—



each of these three variables. Whittaker (1990) determines which variables influence  $a$  by considering the log-linear model  $ac|ad|ae|ag$  induced by the generators of model (5.1) that involve  $a$ . Using maximum likelihood estimation of log-linear parameters, Whittaker obtains the following estimates of the logistic regression of  $a$  on  $c, d, e$  and  $g$ :

$$(5.2) \quad \log \frac{p(a = 1|c, d, e, g)}{p(a = 0|c, d, e, g)} = \text{const.} - 1.33c - 1.32d + 0.69e - 2.17g.$$

Equation (5.2) seems to support our findings based on CGGMs, as it indicates a negative association between  $(a, c)$ ,  $(a, d)$ ,  $(a, g)$ , and a positive association between  $(a, e)$ . Moreover, the association between  $a$  and  $e$  is the weakest of the four. The CGGMs estimate  $\tilde{\rho}_{a,e} = 0.04$  which is about half of  $\tilde{\rho}_{a,c}$  or  $\tilde{\rho}_{a,d}$ . The absolute values of the regression coefficients in (5.2) share the same pattern.

We remark that Table 3 reports a posterior inclusion probability equal to 0.93 for the edge  $(a, b)$ . However, the CGGMs estimate the pairwise correlation  $\Psi_{a,b}$  to be 0.15 and the Cramér’s V association  $\rho_{a,b}$  to be 0.01. Therefore, the CGGMs do not seem to indicate a relevant interaction between variables  $a$  and  $b$  which is in line with Whittaker’s findings who did not include an interaction term  $ab$  in model (5.1). This represents an example where an edge vanishes as we move from the latent variables space to the observed variables space. We would expect the opposite to happen in most applications, that is, edges or associations could be lost when moving from the observed to the latent variables.

5.2. *The NLTCs functional disability data.* We come back to the  $2^{16}$  functional disability table introduced in Section 1. Dobra, Erosheva and Fienberg (2003) analyze these data from a disclosure limitation perspective, while Fienberg et al. (2010) develop latent class (LC) models that are very similar to the Grade of Membership (GoM) models of Erosheva, Fienberg and Joutard (2007). The need to consider alternatives to log-linear models for the NLTCs data comes from the severe imbalance that exists among the cell counts in this table. The largest cell count is 3853, but most of the cells (62,384 or 95.19%) contain counts of zero, while 1729 (2.64%) contain counts of 1 and 1499 (0.76%) contain counts of 2. There are 24 cells with counts larger than 100, which accounts for 42% of the observed sample size 21,574. This gives a very small mean number of observations per cell of 0.33, which is indicative of an extremely high degree of sparsity that is characteristic of high-dimensional categorical data.

We ran 100 replicates of the Markov chain Monte Carlo sampler from Section 4.1 for 500,000 iterations with a burn-in time of 50,000 iterations. Figure 2 shows the convergence of these Markov chains to the joint distribution (4.5). The mean number of edges of the sampled graphs is 72 or 60% of the total number of edges. Table 7 compares the expected cell values of the six largest counts as estimated with the Grade of Membership (GoM) model of Erosheva, Fienberg and Joutard (2007), the latent class (LC) model of Fienberg et al. (2010) and the

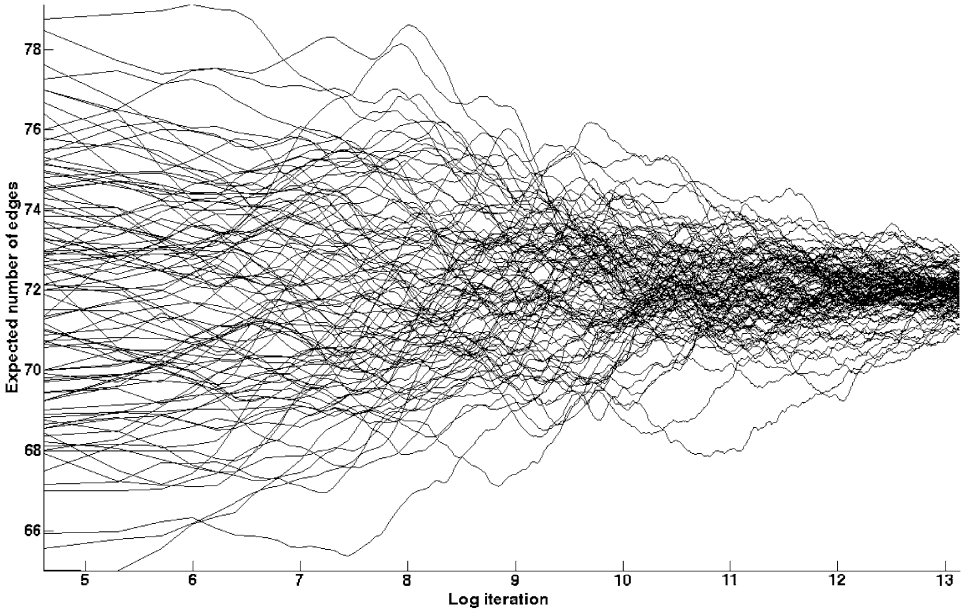


FIG. 2. Estimates of the posterior expected number of edges in the CGGMs for the NLTCs functional disability data.

CGGMs. All three models seem to perform comparably well in terms of capturing the underlying dependency patterns that lead to the largest counts in this  $2^{16}$  table.

In Table 8 we show the association structure of the latent variables  $Z$ . We give posterior estimates of the pairwise correlations  $\Upsilon_{v_1, v_2}$  and posterior inclusion probabilities for each edge  $(v_1, v_2)$ . All the estimates of the pairwise correlations are quite large and strictly positive, which is intuitively correct: the ability to perform any activity of daily living is positively correlated with the ability to perform any other activity. In Table 9 we show the association structure of the observed vari-

TABLE 7

Expected cell counts for the top six largest counts cells in the NLTCs data. We report the results obtained from the GoM model [Erosheva, Fienberg and Joutard (2007)], the LC model [Fienberg et al. (2010)] and the CGGMs. Here 1 stands for healthy and 2 stands for disabled

Cell	Observed	GoM	LC	CGGMs
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	3853	3269	3836.01	3767.76
1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1	1107	1010	1111.51	1145.86
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	660	612	646.39	574.76
1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1	351	331	360.52	452.75
1 1 1 1 1 1 2 1 1 1 1 2 1 1 1 1	303	273	285.27	350.24
1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1	216	202	220.47	202.12

TABLE 8  
*Estimated correlations (elements under the main diagonal) and posterior inclusion probabilities of edges (elements above the main diagonal) in the NLTCs data*

	ADL						IADL									
	1	2	3	4	5	6	1	2	3	4	5	6	7	8	9	10
ADL																
1	—	1	1	0.24	0.46	0.42	1	0.68	0.87	0.98	0.33	1	0.23	0.46	0.15	1
2	0.72	—	1	0.19	0.42	0.94	1	0.10	0.18	0.09	0.10	0.76	1	0.23	0.21	0.17
3	0.78	0.74	—	1	0.36	1	1	0.13	0.50	0.77	0.10	0.74	0.13	0.78	0.24	0.16
4	0.51	0.54	0.64	—	1	1	0.28	0.16	1	1	0.12	0.20	0.36	1	0.14	0.77
5	0.33	0.43	0.41	0.66	—	0.44	0.15	0.54	0.30	0.95	0.18	1	0.81	1	0.96	1
6	0.62	0.65	0.73	0.82	0.66	—	1	1	0.34	0.82	0.10	0.63	0.93	0.81	0.16	0.27
IADL																
1	0.74	0.77	0.76	0.68	0.58	0.83	—	1	1	0.67	0.19	0.20	0.19	0.21	0.95	1
2	0.64	0.69	0.68	0.68	0.62	0.82	0.88	—	1	0.30	0.13	0.19	0.19	0.27	0.55	0.72
3	0.65	0.71	0.66	0.62	0.61	0.79	0.90	0.90	—	1	0.16	0.23	1	0.31	0.92	0.27
4	0.49	0.58	0.55	0.66	0.64	0.76	0.78	0.83	0.87	—	0.12	1	0.42	0.74	0.23	0.44
5	0.45	0.56	0.48	0.52	0.65	0.60	0.65	0.63	0.67	0.61	—	1	1	1	0.97	0.65
6	0.45	0.59	0.52	0.56	0.64	0.64	0.68	0.66	0.70	0.66	0.79	—	1	0.16	0.11	1
7	0.60	0.70	0.60	0.54	0.57	0.65	0.76	0.71	0.77	0.66	0.79	0.79	—	0.33	1	1
8	0.39	0.50	0.43	0.56	0.87	0.63	0.62	0.64	0.66	0.64	0.77	0.72	0.71	—	1	0.84
9	0.48	0.57	0.49	0.55	0.74	0.64	0.67	0.68	0.71	0.65	0.79	0.75	0.80	0.89	—	1
10	0.65	0.69	0.63	0.54	0.52	0.65	0.77	0.70	0.75	0.63	0.74	0.75	0.87	0.68	0.77	—

TABLE 9  
*Estimated Cramér's V associations (elements under the main diagonal) and posterior probabilities  $p(H_{1,\rho}|x^{(1:n)})$  (elements above the main diagonal) in the NLTCs data*

	ADL						IADL									
	1	2	3	4	5	6	1	2	3	4	5	6	7	8	9	10
ADL																
1	—	1	1	0	0	0.61	1	1	1	0	0	0	0.99	0	0	1
2	0.21	—	1	0.43	0	1	1	1	1	0	0.99	1	1	0.08	1	1
3	0.26	0.25	—	1	0	1	1	1	1	0	0.05	0.45	1	0	0.14	0.99
4	0.07	0.10	0.15	—	1	1	1	1	1	1	0.38	1	0.32	1	0.98	0
5	0.03	0.06	0.05	0.21	—	1	0.99	1	0.98	1	1	1	0.32	1	1	0
6	0.10	0.14	0.20	0.38	0.21	—	1	1	1	1	1	1	1	1	1	0.03
IADL																
1	0.21	0.28	0.28	0.18	0.11	0.28	—	1	1	1	1	1	1	1	1	1
2	0.14	0.19	0.19	0.21	0.16	0.34	0.43	—	1	1	1	1	1	1	1	1
3	0.16	0.22	0.18	0.13	0.11	0.22	0.48	0.40	—	1	1	1	1	1	1	1
4	0.04	0.08	0.08	0.19	0.18	0.25	0.15	0.23	0.13	—	0.33	1	0.13	1	1	0
5	0.06	0.12	0.09	0.10	0.14	0.14	0.18	0.17	0.19	0.10	—	1	1	1	1	1
6	0.06	0.12	0.10	0.14	0.19	0.20	0.19	0.21	0.18	0.17	0.27	—	1	1	1	1
7	0.13	0.21	0.14	0.10	0.10	0.14	0.27	0.21	0.28	0.09	0.30	0.23	—	1	1	1
8	0.05	0.09	0.07	0.14	0.39	0.19	0.16	0.19	0.17	0.15	0.26	0.26	0.20	—	1	0.95
9	0.07	0.13	0.09	0.11	0.20	0.16	0.20	0.21	0.23	0.12	0.31	0.25	0.30	0.42	—	1
10	0.14	0.16	0.13	0.06	0.05	0.09	0.20	0.13	0.20	0.05	0.19	0.12	0.32	0.11	0.20	—

ables  $X$ . For every pair  $X_{v_1}$  and  $X_{v_2}$ , we give the posterior means of  $\rho_{v_1, v_2}$  and estimates of the posterior probabilities  $p(H_{1, \rho}^{v_1, v_2} | x^{(1:n)})$  with  $H_{1, \rho}^{v_1, v_2} : \rho_{v_1, v_2} > 0.1$ . The Cramér’s V values indicate that independence is unlikely to hold for any pair of observed variables, which is consistent with the large positive correlations we estimated in the latent space. In fact, 88 pairs of observed variables have a Bayes factor  $B_{\rho}^{v_1, v_2}$  greater than 100, which constitutes strong evidence in favor of the hypothesis  $H_{1, \rho}^{v_1, v_2}$  [Kass and Raftery (1995)]. Thus, the NLTCs data shows that approximately 73% pairs of ADLs and IADLs are certainly not independent of each other.

The topology of the sampled graphs is indicative of the relative importance of each disability measure with respect to the others in the latent variables space. The structure of a graph can be summarized by the number of neighbors of each vertex, that is, the number of edges that involve each variable. This is usually called the degree of a vertex. A larger degree indicates an increased number of interactions in which a latent variable participates. Since in the NLTCs data all the latent variables are positively associated with each other, having one disability increases the likelihood of having other disabilities. The degree of a variable reflects the number of disabilities that are not conditionally independent of this variable given the others.

In the observed variables space we quantify the relative importance of a variable  $X_{v_1}$  as the sum of the Cramér’s V associations  $\rho_{v_1, v_2}$  between  $X_{v_1}$  and some other variable  $X_{v_2}$ . When computing these cumulative Cramér’s V associations we assume that the  $120 - 88 = 32$  pairwise associations with a Bayes factor below 100 are set to zero. Figure 3 shows the posterior expected degrees of the 16 disability

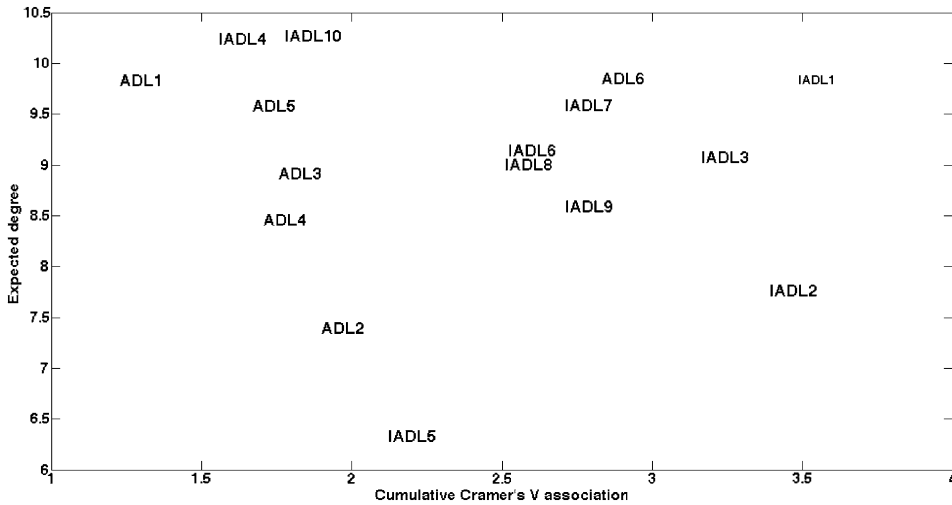


FIG. 3. Cumulative Cramér’s V associations (x-axis) and posterior expected degrees (y-axis) of the 16 disability measures from the NLTCs functional disability data.

measures plotted against the corresponding cumulative Cramér's  $V$  associations. We see that IADL4 (cooking) and IADL10 (telephoning) stand out in the latent space. Most individuals included in the survey (67.6%) are unable to cook, hence, there is no surprise that IADL4 is the second most connected variable. However, only a relatively small number of people (10.6%) cannot use the telephone on their own. In fact, more people are disabled with respect to any of the other 15 measures. As such, it might be counterintuitive to see that IADL10 has the highest degree of connectivity. In the observed variables space the top three cumulative Cramér's  $V$  associations are obtained for IADL1, IADL2 and IADL3. We note that IADL1 (doing heavy house work) and IADL2 (doing light house work) are nested, hence, we would expect their association scores to be related. This indicates a good degree of consistency of the dependency structure identified by the CGGMs. Since IADL1 is also highly connected in the latent space, Figure 3 suggests that IADL1 is key to a principled assessment of the disability level of a person.

The CGGMs clearly show that the 16 disability measures recorded in the NLTCS data should not be treated on an equal footing. Some measures such as IADL1 or IADL10 indicate more serious disabilities than others, which is not necessarily reflected in the number of people reporting that particular disability. Simply counting the number of disabilities a person has can be very misleading when evaluating the overall disability level of an individual. This remark could shed a new light on the findings reported in [Manton and Gu \(2001\)](#) who only make the distinction between ADLs and IADLs.

**6. Discussion.** The inference approach we presented in this paper extends Gaussian graphical models to data sets in which the multivariate normal assumption for the observed variables is unlikely to hold. The CGGMs capture conditional independence relationships among a set of latent variables that are in a one-to-one relationship with the set of observed variables. The fact that the number of latent variables coincides with the number of observed variables avoids the difficult statistical issue of having to select the number of latent classes—see the excellent discussions in [Erosheva, Fienberg and Joutard \(2007\)](#) and [Fienberg et al. \(2010\)](#).

Our goal was to model dependencies separately from the univariate marginal distribution of each variable. As such, we did not include a parametric representation of the marginal distributions in our framework. [Pitt, Chan and Kohn \(2006\)](#) give a Bayesian approach to model conditional independence relationships in Gaussian copulas in which the univariate marginal distributions are allowed to depend on a set of parameters and on certain sets of explanatory variables. There is a definite possibility to combine our prior specification for the precision matrix for the latent variables with the methods of [Pitt, Chan and Kohn \(2006\)](#) into a procedure that takes into account the uncertainty in the specification of the univariate distributions.

The CGGMs are applicable to any observational study for the purpose of identifying conditional independence relationships. The only requirement is that the

observed variables are binary, ordinal or continuous. The extended rank likelihood [Hoff (2007)] is a key component of our framework. A necessary condition for its correct application is that there exists an ordering of the possible values of any observed variable—see Section 3. Our framework does not allow the presence of discrete variables that are not binary or ordinal.

Although the interactions among the latent variables do not go beyond second-order moments, CGGMs give sensible results in the analysis of sparse contingency tables because they allow inference through Bayesian model averaging. By contrast, log-linear models contain higher-order interaction terms but model averaging is no longer an option: the same interaction term has a different interpretation in various log-linear models. As such, one has to choose one log-linear model and perform inference given this single model. When the sample size is small with respect to the total number of possible models, such a determination might not be appropriate. The data might not contain enough information to distinguish between log-linear models that are very close to each other and have almost the same posterior probability—see, for example, the analysis of the Rochdale data from Dobra and Massam (2010). Our use of CGGMs does not involve choosing one particular model, but averaging with respect to many models on the latent space. We hope that CGGMs will play a significant role in many quantitative fields of research.

**Acknowledgments.** The authors thank Peter Hoff for useful discussions. The authors are also grateful to Elena Erosheva who provided the NLTCs data. The authors thank the Editor and anonymous reviewers for their comments that improved the quality of this writing.

#### SUPPLEMENTARY MATERIAL

**Supplement: C++ implementation of copula Gaussian graphical models** (DOI: [10.1214/10-AOAS397SUPP](https://doi.org/10.1214/10-AOAS397SUPP); .zip). We provide source code for the methodology described in this paper. Our program takes advantage of cluster computing to run several Markov chains in parallel. By using this code, one can replicate the analyses of the Rochdale data and the NLTCs functional disability data for which we give sample input files.

#### REFERENCES

- ATAY-KAYIS, A. and MASSAM, H. (2005). A Monte Carlo method for computing the marginal likelihood in nondecomposable Gaussian graphical models. *Biometrika* **92** 317–335. [MR2201362](#)
- CRAMÉR, H. (1946). *Mathematical Methods of Statistics*. Princeton Univ. Press, Princeton, NJ. [MR0016588](#)
- DIACONIS, P. and YLVIKAKER, D. (1979). Conjugate priors for exponential families. *Ann. Statist.* **7** 269–281. [MR0520238](#)
- DOBRA, A., EROSHEVA, E. A. and FIENBERG, S. E. (2003). Disclosure limitation methods based on bounds for large contingency tables with application to disability data. In *Proceedings of Conference on the New Frontiers of Statistical Data Mining* (E. H. Bozdogan, ed.) 93–116. CRC Press, New York. [MR2048950](#)

- DOBRA, A. and LENKOSKI, A. (2010). Supplement to “Copula Gaussian graphical models and their application to modeling functional disability data.” DOI: [10.1214/10-AOAS397SUPP](https://doi.org/10.1214/10-AOAS397SUPP).
- DOBRA, A. and MASSAM, H. (2010). The mode oriented stochastic search algorithm (MOSS) for log-linear models with conjugate priors. *Statist. Methodol.* **7** 240–253.
- DUNSON, D. B. (2006). Bayesian dynamic modeling of latent trait distributions. *Biostatistics* **7** 551–568.
- DUNSON, D. B. and XING, C. (2009). Nonparametric Bayes modeling of multivariate categorical data. *J. Amer. Statist. Assoc.* **104** 1042–1051. [MR2562004](#)
- EROSHEVA, E. A., FIENBERG, S. E. and JOUTARD, C. (2007). Describing disability through individual-level mixture models for multivariate binary data. *Ann. Appl. Statist.* **1** 502–537. [MR2415745](#)
- FIENBERG, S. E., HERSH, P., RINALDO, A. and ZHOU, Y. (2010). Maximum likelihood estimation in latent class models for contingency table data. In *Algebraic and Geometric Methods in Statistics* (P. Gibilisco, E. Riccomagno, M. P. Rogantin and E. H. P. Wynn, eds.) 27–62. Cambridge Univ. Press, Cambridge. [MR2642657](#)
- GENEST, C. and NESLEHOVÁ (2007). A primer on copulas for count data. *Astin Bulletin* **37** 475–515. [MR2422797](#)
- GIUDICI, P. and GREEN, P. J. (1999). Decomposable graphical Gaussian model determination. *Biometrika* **86** 785–801. [MR1741977](#)
- GREEN, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82** 711–732. [MR1380810](#)
- HOFF, P. D. (2007). Extending the rank likelihood for semiparametric copula estimation. *Ann. Appl. Statist.* **1** 265–283. [MR2393851](#)
- JONES, B., CARVALHO, C., DOBRA, A., HANS, C., CARTER, C. and WEST, M. (2005). Experiments in stochastic computation for high-dimensional graphical models. *Statist. Sci.* **20** 388–400. [MR2210226](#)
- KASS, R. and RAFTERY, A. E. (1995). Bayes factors. *J. Amer. Statist. Assoc.* **90** 773–795.
- LAURITZEN, S. L. (1996). *Graphical Models*. Oxford Univ. Press, Oxford. [MR1419991](#)
- LENKOSKI, A. and DOBRA, A. (2010). Computational aspects related to inference in Gaussian graphical models with the G-Wishart prior. *J. Comput. Graph. Statist.* DOI: [10.1198/jcgs.2010.08181](https://doi.org/10.1198/jcgs.2010.08181).
- LETAC, G. and MASSAM, H. (2007). Wishart distributions for decomposable graphs. *Ann. Statist.* **35** 1278–1323. [MR2341706](#)
- LIU, H., LAFFERTY, J. and WASSERMAN, L. (2009). The nonparanormal: Semiparametric estimation of high dimensional undirected graphs. *J. Mach. Learn. Res.* **10** 2295–2328. [MR2563983](#)
- MADIGAN, D. and YORK, J. (1995). Bayesian graphical models for discrete data. *Int. Statist. Rev.* **63** 215–232.
- MANTON, K. G., CORDER, L. and STALLARD, E. (1993). Estimates of change in chronic disability and institutional incidence and prevalence rate in the US elderly populations from 1982 to 1989. *J. Gerontol. Soc. Sci.* **48** S153–S166.
- MANTON, K. G. and GU, X. (2001). Changes in prevalence of chronic disability in the United States black and nonblack population above age 65 from 1982 to 1999. *Proc. Natl. Acad. Sci. USA* **98** 6354–6359.
- MUIRHEAD, R. J. (2005). *Aspects of Multivariate Statistical Theory*. Wiley, New York. [MR0652932](#)
- MUTHÉN, B. (1984). A general structural equation model with dichotomous, ordered categorical, and continuous latent variables indicators. *Psychometrika* **49** 115–132.
- NELSEN, R. B. (1999). *An Introduction to Copulas*. Springer, New York. [MR1653203](#)
- PITT, M., CHAN, D. and KOHN, R. (2006). Efficient Bayesian inference for Gaussian copula regression models. *Biometrika* **93** 537–554. [MR2261441](#)



- ROVERATO, A. (2002). Hyper inverse Wishart distribution for non-decomposable graphs and its application to Bayesian inference for Gaussian graphical models. *Scand. J. Statist.* **29** 391–411. [MR1925566](#)
- SCOTT, J. G. and BERGER, J. O. (2006). An exploration of aspects of Bayesian multiple testing. *J. Statist. Plann. Inference* **136** 2144–2162. [MR2235051](#)
- SONG, P. X. K. (2000). Multivariate dispersion models generated from Gaussian copula. *Scand. J. Statist.* **27** 305–320. [MR1777506](#)
- WHITTAKER, J. (1990). *Graphical Models in Applied Multivariate Statistics*. Wiley, New York. [MR1112133](#)
- WONG, F., CARTER, C. K. and KOHN, R. (2003). Efficient estimation of covariance selection models. *Biometrika* **90** 809–830. [MR2024759](#)

DEPARTMENT OF STATISTICS  
DEPARTMENT OF BIOBEHAVIORAL NURSING  
AND HEALTH SYSTEMS  
AND  
CENTER FOR STATISTICS  
AND THE SOCIAL SCIENCES  
UNIVERSITY OF WASHINGTON  
Box 354322  
C-14B PADEL FORD HALL  
SEATTLE, WASHINGTON 98195-4322  
USA  
E-MAIL: [adobra@uw.edu](mailto:adobra@uw.edu)  
URL: <http://www.stat.washington.edu/adobra>

DEPARTMENT OF APPLIED MATHEMATICS  
HEIDELBERG UNIVERSITY  
IM NEUENHEIMER FELD 294  
69120 HEIDELBERG  
GERMANY  
E-MAIL: [lenkoski@stat.washington.edu](mailto:lenkoski@stat.washington.edu)