

STRUCTURE ESTIMATION FOR MIXED GRAPHICAL MODELS IN HIGH-DIMENSIONAL DATA

BY JONAS M. B. HASLBECK

Utrecht University

AND

BY LOURENS J. WALDORP

University of Amsterdam

Undirected graphical models are a key component in the analysis of complex observational data in a large variety of disciplines. In many of these applications one is interested in estimating the undirected graphical model underlying a distribution over variables with different domains. Despite the pervasive need for such an estimation method, to date there is no such method that models all variables on their proper domain. We close this methodological gap by combining a new class of mixed graphical models with a structure estimation approach based on generalized covariance matrices. We report the performance of our methods using simulations, illustrate the method with a dataset on Autism Spectrum Disorder (ASD) and provide an implementation as an R-package.

1. Introduction. Determining conditional independence relationships through undirected graphical models, also known as Markov random fields (MRFs), is a key component of the statistical analysis of complex observational data in a wide variety of domains such as statistical physics, image analysis, medicine and more recently psychology (Borsboom and Cramer, 2013). In many of these applications one is interested in estimating the MRF underlying a joint distribution over *variables with different domains*.

As an example, consider a dataset of questionnaire responses of individuals diagnosed with Autism Spectrum Disorder (ASD), covering demographics, social environment, diagnostic measurements and aspects of well-being. A central research question in the study of Autism is how to explain individual differences in well-being. Many studies tried to answer this question by focusing on the relation between well-being and specific areas such as social functioning, cognitive ability, education and working conditions (e.g. Magiati, Tay and Howlin, 2014; Anderson, Liang and Lord, 2014). While

MSC 2010 subject classifications: Primary Structure estimation, mixed distributions, graphical models; secondary Generalized covariance matrices

this approach provides valuable insights, one necessarily misses the full picture of all variables and might misinterpret relationships that change when additional variables are taken into account. An alternative is an integrated analysis, in which one determines the conditional independence relationships of *all variables* by estimating the MRF underlying the multivariate distribution over all variables. In most datasets, however, this means that we need a method to estimate an MRF underlying a multivariate distribution over variables of different domains. In our example dataset, for instance, we have variables that are continuous (age), ordinal (IQ-scores), categorical (type of housing) and count-valued (number of different medications).

Despite the pervasive need for a method to estimate a MRF underlying a joint distribution over mixed variables in many disciplines, so far such a general method is not available. In the present paper we address this methodological gap by combining a new class of *mixed joint distributions* (see Section 2.2) with a structure estimation approach based on *generalized covariance matrices* (see Section 2.3). We thereby provide the first method that estimates a mixed MRF in which all variables are modeled on their proper domain. This avoids possible information loss due to variable transformations. In addition, our method is attractive in terms of interpretation as it is similar to the well-known Gaussian case. We introduce basic concepts in Section 2.1, describe our estimation algorithm in detail in Section 2.5, show performance benchmarks in Section 3 and apply our method to the above dataset on ASD in Section 4. In the remainder of this section, we review existing methods to estimate MRFs underlying mixed distributions.

Graphical models associated with distributions over one type of variable are well-known and used in many applications. The most prominent example is the Gaussian MRF underlying the multivariate Gaussian distribution. A corollary of the Hammersley-Clifford theorem (Lauritzen, 1996) states that zeros in the inverse covariance matrix of the multivariate Gaussian distribution indicate absent edges in the corresponding graphical model. Two classes of efficient algorithms leverage this relationship in order to estimate the Gaussian MRF: *global* methods estimate the whole graph by directly estimating the inverse covariance matrix using a penalized likelihood (Yuan and Lin, 2007; Banerjee, El Ghaoui and d’Aspremont, 2008; Friedman, Hastie and Tibshirani, 2008) and *nodewise* methods estimate the neighborhood of each node separately by solving a collection of sparse regression problems using the Lasso (Meinshausen and Bühlmann, 2006). Another graphical model that is widely used is the MRF underlying the Potts model (see e.g. Wainwright and Jordan, 2008). For the Ising model (Potts model with $m = 2$ categories), Ravikumar, Wainwright and Lafferty (2010) proposed a node-

wise estimation method based on ℓ_1 -regularized neighborhood regression. This method was further extended by [van Borkulo et al. \(2014\)](#), who select the regularization parameter using the Extended Bayesian Information Criterion (EBIC), which is known to perform well in selecting sparse graphs ([Foygel and Drton, 2014](#)). For general multivariate discrete distributions, [Loh and Wainwright \(2013\)](#) proposed an approach based on the estimation of the inverses of generalized covariance matrices.

Work on mixed graphical models is rather recent and many of the existing methods are based on non-parametric extensions of the graphical models described above: the non-paranormal ([Liu, Lafferty and Wasserman, 2009](#); [Lafferty, Liu and Wasserman, 2012](#)) method uses transforms that Gaussianize the data and then fits a Gaussian MRF and thereby offers a graphical model consisting of different distributions with continuous domain. Similarly, the copula-based method of [Dobra and Lenkoski \(2011\)](#) uses thresholded latent Gaussian variables to model ordinal variables. Other methods use non-parametric approximations such as rank-based estimators to the correlation matrix, and then fit a Gaussian MRF ([Xue and Zou, 2012](#); [Liu et al., 2012](#)).

Another way of modeling mixed graphical models is to sidestep multivariate densities altogether and relate a set of multivariate response variables of one type to multivariate covariate variables of another type. This can be done by using multiple regression or multi-task learning models ([Evgeniou and Pontil, 2007](#)). More recent approaches also allow these multiple regression models to associate covariates with mixed types of responses ([Yang, Kim and Xing, 2009](#)). Also, in many machine learning learning procedures, mixed types of variables are accounted for implicitly by using suitable distance- or entropy-based measures ([Hastie, Tibshirani and Friedman, 2009](#); [Hsu, Chen and Su, 2007](#)). However, the sample complexity of non-parametric methods is typically inferior to those that learn parametric models, especially in high-dimensional settings.

Parametric approaches involve methods based on latent variables that permit mixed continuous and count variables ([Sammel, Ryan and Legler, 1997](#)) or dependencies between exponential family members through a latent Gaussian MRF ([Rue, Martino and Chopin, 2009](#)). While these approaches provide statistical models for mixed data, they model dependencies between observed variables using a latent layer of unobserved variables; this typically renders techniques computationally expensive and possibly intractable when estimating these models with strong statistical guarantees. A classic model for mixed data *without* latent variables is the conditional Gaussian model, that combines categorical and Gaussian variables ([Lauritzen, 1996](#)). Here,

Gaussian variables are modeled as a multivariate Gaussian that is conditioned on all possible states of the categorical variable, which renders the model computationally intractable. The model becomes more feasible when it is restricted to pairwise or three-way interactions (Lee and Hastie, 2012; Cheng, Levina and Zhu, 2013, respectively); however, the general problem remains. Finally, Yang et al. (2014) introduces a way to combine different conditional distributions to their joint distribution, given that each conditional is an exponential family member.

The key idea of our paper is to generalize the *generalized covariance approach* proposed by Loh and Wainwright (2013) for discrete random variables to the class of *mixed joint distributions* over random variables from the exponential family introduced by Yang et al. (2014). We thereby obtain an estimation algorithm for mixed graphical models that both models all variables on their proper domain and is easy to interpret as the estimation method is similar to the well-known Gaussian case.

2. Estimation of mixed Markov Random Fields in high-dimensional data.

2.1. *Markov random fields.* Undirected graphical models or Markov random fields (MRFs) are families of probability distributions that respect the structure of an undirected graph. An undirected graph $G = (V, E)$ consists of a collection of nodes $V = \{1, 2, \dots, p\}$ and a collection of edges $E \subseteq V \times V$. A *node cutset* is a subset U of nodes that breaks the graph into two or more nonempty components when it is removed from the graph. A *clique* is a subset $C \subseteq V$ such that $(s, t) \in E$ for all $s, t \in C$ where $s \neq t$. A *maximal clique* is a clique that is not properly contained within any other clique. The neighborhood $N(s)$ of node s is defined as the set of nodes that are connected to s by an edge, $N(s) := \{t \in V \mid (s, t) \in E\}$. The degree of a node s is denoted by $\text{deg}(s) = |N(s)|$. Throughout the paper we use the shorthand $X_{\setminus s}$ for $X_{V \setminus \{s\}}$.

To each vertex s in graph G we associate a random variable X_s taking values in a space \mathcal{X} . For any subset $A \subseteq V$, we use the shorthand $X_A := \{X_s, s \in A\}$. For three subsets of nodes, A , B , and U , we write $X_A \perp\!\!\!\perp X_B \mid X_U$ to indicate that the random vector X_A is independent of X_B when conditioning on X_U . Markov random fields can be defined in terms of the global Markov property:

DEFINITION 1. (*Global Markov property*). If $X_A \perp\!\!\!\perp X_B \mid X_U$ whenever U is a vertex cutset that breaks the graph into disjoint subsets A and B , then

the random vector $X := (X_1, \dots, X_p)$ is Markov with respect to the graph G .

Note that the neighborhood set $N(s)$ is always a vertex cutset for $A = \{s\}$ and $B = V \setminus \{s \cup N(s)\}$.

By the Hammersley-Clifford Theorem (e.g. [Lauritzen, 1996](#)), for strictly positive probability distributions, the global Markov property is equivalent to the Markov factorization property. Consider for each clique $C \in \mathcal{C}$ a clique-compatibility function $\psi_C(X_C)$ that maps configurations $x_C = \{x_s, s \in V\}$ to \mathbb{R}^+ such that ψ_C only depends on the variables X_C corresponding to the clique C .

DEFINITION 2. (*Markov factorization property*). *The distribution of X factorizes according to G if it can be represented as a product of clique functions*

$$(1) \quad P(X) \propto \prod_{C \in \mathcal{C}} \psi_C(X_C).$$

Because we focus on strictly positive distributions, we can represent (1) in terms of an exponential family associated with the cliques C in G

$$(2) \quad P(X) = \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C \phi_C(X_C) - \Phi(\theta) \right\},$$

where the functions $\phi_C(X_C) = \log \psi_C(X_C)$ are sufficient statistic functions specified by the exponential family member at hand, θ_C are parameters associated with these functions and $\Phi(\theta)$ is the log-normalizing constant

$$\Phi(\theta) = \log \int_{\mathcal{X}} \sum_{C \in \mathcal{C}} \theta_C \phi_C(X_C) \nu(dx).$$

2.2. Mixed joint distributions. [Yang et al. \(2014\)](#) introduced a special case of the form (2) which allows to model any combination of conditional univariate members of the exponential family within one joint distribution which respects the neighborhoods of the conditional distributions. We first describe this class of models and then provide an example.

Consider a p -dimensional random vector $X = (X_1, \dots, X_p)$ with each variable X_s taking values in a potentially different set \mathcal{X}_s and let $G = (V, E)$ be an undirected graph over p nodes corresponding to the p variables. Now

suppose the node-conditional distribution of node X_s conditioned on all other variables $X_{\setminus s}$ is given by an arbitrary univariate exponential family distribution

$$(3) \quad P(X_s|X_{\setminus s}) = \exp \{E_s(X_{\setminus s})\phi_s(X_s) + C_s(X_s) - \Phi(X_{\setminus s})\},$$

where the functions of the sufficient statistic $\phi_s(\cdot)$ and the base measure $C_s(\cdot)$ are specified by the choice of exponential family and the canonical parameter $E_s(X_{\setminus s})$ is a function of all variables except X_s .

These node-conditional distributions are consistent with a joint MRF distribution over the random vector X as in (1), that is, Markov with respect to graph $G = (V, E)$ with clique-set \mathcal{C} of size at most k , if and only if the canonical parameters $\{E_s(\cdot)\}_{s \in V}$ are a linear combination of k -th order products of univariate sufficient statistic functions $\{\phi(X_t)\}_{t \in N(s)}$

$$\theta_s + \sum_{t \in N(s)} \theta_{st} \phi_t(X_t) + \dots + \sum_{t_1, \dots, t_{k-1} \in N(s)} \theta_{t_1, \dots, t_{k-1}} \prod_{j=1}^{k-1} \phi_{t_j}(X_{t_j}),$$

where $\theta_s := \{\theta_s, \theta_{st}, \dots, \theta_{st_2 \dots t_k}\}$ is a set of parameters and $N(s)$ is the set of neighbors of node s according to graph G . The corresponding joint distribution has the form

$$(4) \quad P(X; \theta) = \exp \left\{ \sum_{s \in V} \theta_s \phi_s(X_s) + \sum_{s \in V} \sum_{t \in N(s)} \theta_{st} \phi_s(X_s) \phi_t(X_t) + \dots + \sum_{t_1, \dots, t_k \in \mathcal{C}} \theta_{t_1, \dots, t_k} \prod_{j=1}^k \phi_{t_j}(X_{t_j}) + \sum_{s \in V} C_s(X_s) - \Phi(\theta) \right\},$$

where $\Phi(\theta)$ is the log-normalization constant.

A limitation of this class of models is that in case they include two univariate distributions with infinite domain, they are not normalizable if *neither* both distributions are infinite only from one side *nor* the base measures are bounded with respect to the moments of the random variables (for details see [Yang et al., 2014](#)). We return to this limitation in the discussion.

As a concrete example of (4), take the Ising-Gaussian model: consider a random vector $X := (Y, Z)$, where $Y = \{Y_1, \dots, Y_p\}$ are univariate Gaussian random variables, $Z = \{Z_1, \dots, Z_p\}$ are univariate Bernoulli random

variables and we only consider pairwise interactions between sufficient statistics. For the univariate Gaussian distribution (with known σ^2) the sufficient statistic function is $\phi_Y(Y_s) = \frac{Y_s}{\sigma_s}$ and the base measure is $C_Y(Y_s) = -\frac{Y_s^2}{2\sigma_s^2}$. The Bernoulli distribution has the sufficient statistic function $\phi_{Z_r} = Z_r$ and the base measure $C_Z(Z_r) = 0$. From (4) follows that this mixed distribution has the form

$$(5) \quad P(Y, Z) \propto \exp \left\{ \sum_{s \in V_Y} \frac{\theta_s}{\sigma_s} Y_s + \sum_{r \in V_Z} \theta_r Z_r + \sum_{(s,t) \in E_Y} \frac{\theta_{st}}{\sigma_s \sigma_t} Y_s Y_t + \sum_{(r,q) \in E_Z} \theta_{rq} Z_r Z_q + \sum_{(s,r) \in E_{YZ}} \frac{\theta_{sr}}{\sigma_s} Y_s Z_r - \sum_{s \in V_Y} \frac{Y_s^2}{2\sigma_s^2} \right\}.$$

If X_r is a Bernoulli random variable, the node-conditional has the form

$$P(X_r | X_{\setminus r}) \propto \exp \left\{ \theta_r Z_r + \sum_{q \in N(r)_Z} \theta_{rq} Z_r Z_q + \sum_{t \in N(r)_Y} \frac{\theta_{rt}}{\sigma_t} Z_r Y_t \right\}.$$

Note that this form is equivalent to the node-conditional Ising model with one term added for interactions between Bernoulli and Gaussian random variables.

If X_s is a Gaussian random variable, the node-conditional has the form

$$P(X_s | X_{\setminus s}) \propto \exp \left\{ \frac{\theta_s}{\sigma_s} Y_s + \sum_{t \in N(s)_Y} \frac{\theta_{st}}{\sigma_s \sigma_t} Y_s Y_t + \sum_{r \in N(s)_Z} \frac{\theta_{sr}}{\sigma_s} Y_s Z_r - \frac{Y_s^2}{2\sigma_s^2} \right\}.$$

Not, let $\sigma = 1$, factor out Y_s and let $\mu_s = \theta_s + \sum_{t \in N(s)_Y} \theta_{st} Y_t + \sum_{r \in N(s)_Z} \theta_{sr} Z_r$. Finally, when taking $\frac{\mu_s^2}{2}$ out of the log normalization constant, we arrive with basic algebra at the well-known form of the univariate Gaussian distribution with unit variance

$$P(X_s | X_{\setminus s}) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(X_s - \mu_s)^2}{2} \right\}.$$

In order to estimate the Markov random field underlying these mixed joint distributions, we use generalized covariance matrices. We introduce these in the following section.

2.3. *Generalized covariance matrices of mixed joint distributions.* Consider a random vector $\Psi = \{X_1, \dots, X_p\}$ that is jointly Gaussian. Then the inverse Γ of the covariance matrix $\text{cov}(\Psi)$ is graph-structured in the sense that $\Gamma_{st} = 0$ whenever $(s, t) \notin E$ (e.g. [Lauritzen, 1996](#)). While this result is well-known and leveraged by many structure estimation algorithms for the Gaussian case, the relationship between inverse covariance and graph-structure is generally unknown for general multivariate distributions.

[Loh and Wainwright \(2013\)](#) improved this situation by showing that the inverses of covariance matrices over discrete random vectors are also graph-structured, when augmenting the covariance matrix appropriately with higher order interactions. In the present paper, we extend this result to the more general class of mixed joint distributions as in (4). This will allow us to use the nodewise graph algorithm proposed by [Loh and Wainwright \(2013\)](#) also for the estimation of mixed Markov random fields underlying mixed distributions. In the remainder of this section we first introduce necessary concepts, then illustrate the method using an example and finally state all results formally.

We require the notions of *triangulation*, *junction trees* and *block graph structure*. Triangulation can be defined in terms of chordless cycles, which are sequences of distinct nodes $\{s_1, \dots, s_\ell\}$ such that $(s_i, s_{i+1}) \in E$ for all $1 \leq i \leq \ell - 1$, $(s_\ell, s_1) \in E$ and no other nodes in the cycle are connected by an edge. Given an undirected graph $G = (V, E)$, a triangulation is an augmented graph $\tilde{G} = (V, \tilde{E})$ that does not contain chordless cycles with length larger than 3.

Any triangulation \tilde{G} gives rise to a junction tree representation of G . Nodes in the junction tree are subsets of V corresponding to maximal cliques in \tilde{G} . The intersection of two adjacent cliques C_1 and C_2 is called *separator set* $S = C_1 \cap C_2$. Second, any junction tree must satisfy the *running intersection property*, that is, for each pair U, V of cliques with intersection S , all cliques on the path between U and V contain S . We refer the reader to [Koller and Friedman \(2009\)](#) and [Cowell et al. \(2007\)](#) for a more detailed treatment of the notions of triangulation and junction trees.

Consider a random vector $\Psi(X)$ that factors according to a graph G and let \mathcal{C} be the set of *all* cliques in G . We say that an inverse covariance matrix $\Gamma = (\text{cov}(\Psi(X))^{-1})$ is *block graph-structured*, if both of the following statements are true for any $A, B \in \mathcal{C}$:

1. If A, B are not subsets of the same maximal clique, the entry $\Gamma(A, B)$ is identically zero.
2. For almost all parameters θ , the entry $\Gamma(A, B)$ is nonzero whenever A and B belong to the same maximal clique.

Consider an Ising-Gaussian model as in (5) that factors according to the graph in Figure 1 (a), where X_1, X_3 are Bernoulli random variables and X_2, X_4 are Gaussian random variables. We let all node parameters be $\theta_s = 0.1$ for all $s \in V$ and the edge parameters be $\theta_{st} = 0.5$ for all $(s, t) \in E$. Recall that if X was distributed jointly Gaussian, standard theory would predict that Γ is graph-structured. The empirically computed inverse covariance matrix Γ of this model is displayed in Figure 1 (b). While the absent edge (4,2) between the Gaussian nodes is reflected by a zero entry at $\Gamma_{4,1}$ this is not the case for the other absent edge (3,1) and Γ is therefore not graph-structured.

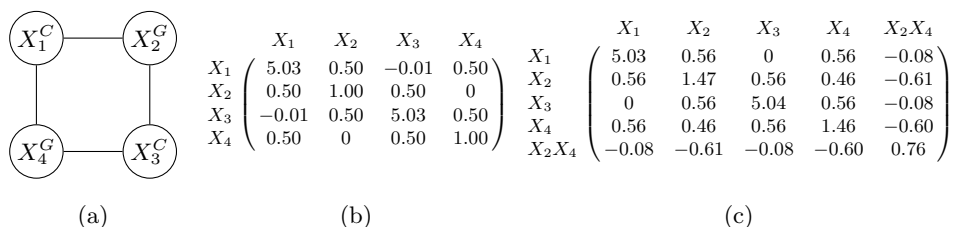


FIG 1. (a) Ising-Gaussian Markov random field, (b) Inverse covariance matrix, (c) Inverse of augmented covariance matrix.

Now we consider the same covariance matrix but augment it with the interaction X_2X_4 and compute its inverse Γ , which is displayed in 1 (c). We now have a zero entry at $\Gamma_{3,1}$, reflecting the absent edge at (3,1). However, there is a nonzero entry at $\Gamma_{4,2}$, which does not reflect the absent edge (4,2). We see that augmenting interactions to the covariance matrix renders its inverse graph structured with respect to the new, triangulated graph \tilde{G} .

We will now state the general result regarding the relationship between the inverses of augmented covariance matrices and the underlying graph structure. Equipped with this result we will come back to the above example and explain the zero-pattern in 1(c).

Let $\mathcal{A} \subseteq \mathcal{C}$ be a set of cliques and define the random vector $\Psi(\phi(X); \mathcal{A}) := \{\phi_C(X_C), C \in \mathcal{A}\}$.

THEOREM 1. *Consider a mixed joint distribution $P_\theta(X)$ as in (4) that factorizes according to a triangulated graph \tilde{G} and let $\tilde{\mathcal{C}}$ be the set of all cliques in \tilde{G} . Then the generalized covariance matrix $\text{cov}(\Psi(\phi(X); \tilde{\mathcal{C}}))$ is invertible, and its inverse $\tilde{\Gamma}$ is block graph structured.*

The proof is based on an equality between the inverse covariance $\tilde{\Gamma}$ and the conjugate dual $\Phi^*(\mu)$ of the log-normalizing constant $\Phi(\theta)$. We impose

the triangulation condition, because it is a sufficient condition for being able to verify the claims of Theorem 1 in $\Phi^*(\mu)$. We provide a proof for Theorem 1 in Section 2.4.1.

Returning to the example in Figure 1, if we take the triangulation \tilde{G} in which we augment the edge (4,2), and let $\phi(X) = X$, then the set of all cliques in \tilde{G} is equal to

$$\tilde{\mathcal{C}} = \{X_1, X_2, X_3, X_4, X_1X_2, X_2X_3, X_3X_4, X_4X_1, X_4X_2, X_1X_2X_4, X_2X_3X_4\}.$$

It can now be shown empirically that the 11×11 inverse covariance matrix $\tilde{\Gamma} = (\text{cov}(\Psi(\phi(X); \tilde{\mathcal{C}})))^{-1}$ reflects the graph structure of \tilde{G} : there are zeros at the positions $\Gamma_{A,B}$, corresponding to the functions $\phi_A(X_A) = \prod_{s \in A} \phi_s(X_s)$ and $\phi_B(X_B) = \prod_{s \in B} \phi_s(X_s)$, whenever A and B are not contained in the same maximal clique. This would be the case e.g. for $A = \{1\}$ and $B = \{3\}$.

Theorem 1 requires that we augment the covariance matrix with all cliques $C \in \tilde{\mathcal{C}}$. A corollary of Theorem 1 (see Corollary A.2) states that it suffices to add all non-empty subsets of all separator sets S in \tilde{G} to render Γ graph structured. Stated formally, let $\mathbf{pow}(\mathcal{A}) = \bigcup_{C \in \mathcal{A}} \mathbf{pow}(C)$ be the union of all $2^{|C|} - 1$ nonempty subsets of all cliques $C \in \mathcal{A}$. If \mathcal{S} is the set of all separator sets in \tilde{G} , then we have to augment the original covariance matrix with sufficient statistics associated with the cliques in $\mathbf{pow}(\mathcal{S})$ to render its inverse graph-structured.

Note that we satisfy this condition in the example in Figure 1: We triangulate G by augmenting the edge (3,1) and compute the inverse $\tilde{\Gamma}$ of the covariance matrix over the augmented random vector $\Psi\{\phi(X); V \cup \mathbf{pow}(\mathcal{A})\} = \Psi\{\phi_1(X_1), \phi_2(X_2), \phi_3(X_3), \phi_4(X_4), \phi_4(X_4)\phi_2(X_2)\}$. As we have $(3,1) \notin \tilde{E}$ we predict $\Gamma_{3,1} = 0$, which is what we obtain empirically in Figure 1 (c).

So far we have seen that Theorem 1 enables us to construct a covariance matrix whose inverse reflects the graph-structure of a triangulated graph \tilde{G} . However, in practice, we are not interested in the graph structure of a triangulation of G , but in the graph structure of the *original graph* G . Using a neighborhood selection approach, the following corollary of Theorem 1 enables us to estimate the *original graph* G . Let $\mathcal{A}(s; d) := \{U \subseteq V \setminus \{s\}, |U| = d\}$, the set of all candidate neighborhoods of node s with size equal to d .

COROLLARY 1. *For any node $s \in V$ with $\text{deg}(s) \leq d$ in any graph, the inverse Γ of the covariance matrix $\text{cov}(\Psi(X; \{s\} \cup \mathbf{pow}(\mathcal{A}(s; d))))$ is graph structured such that, $\Gamma(\{s\}, B) = 0$ whenever $\{s\} \neq B \not\subseteq N(s)$. Specifically, $\Gamma(\{s\}, \{t\}) = 0$ for all $t \notin N(s)$.*

This result follows from Theorem 1 by constructing a particular junction tree, in which the set of candidate neighborhoods defined in Corollary 1 separates the node s from the rest of graph G . This new graph consisting of the node s and all cliques $S(s; d)$ is triangulated, because it consists only cliques.

Recall that the candidate neighborhood is required to separate node s from the rest of graph G . This implies that we could define a much smaller set of candidate neighborhoods if we made an additional assumption about the size of the largest clique in graph G that does not include node s . For example, if we know that the true graph reflects the conditional independence structure of a pairwise model, it suffices to define the candidate neighborhood as $\mathbf{pow}(\mathcal{A}(s; 2))$, regardless of $\deg(s)$. We return to this issue in the discussion of Algorithm 1 in Section 2.5.

We can now recover the original graph G by applying Corollary 1 to each node $s \in V$ and take the row s of the inverse Γ as the support of the neighborhood $N(s)$. Because $\Gamma_{s, \setminus s}$ is a scalar multiple of the regression vector of X_s upon $X_{\setminus s}$, we can use linear regression in order to estimate $\Gamma_{s, \setminus s}$. Before providing a detailed description of our nodewise method for graph-estimation in Section 2.5, we provide the proof for Theorem 1 and its corollaries in the following section.

2.4. *Proof of Theorem 1 and Corollary 1.* Theorem 1 is an extension of Theorem 1 in Loh and Wainwright (2013) in that the following theorem is about mixed Markov random fields instead of discrete Markov random fields. We first provide a proof for Theorem 1 and then show that Corollary 1 follows.

2.4.1. *Proof of Theorem 1.* We follow the proof of Loh and Wainwright (2013), which consists of two parts. First, we have establish the equality of the inverse covariance matrix of sufficient statistics $(\text{cov}_\theta\{\phi(X)\})^{-1}$ and the Hessian of the conjugate dual $\nabla^2\Phi^*(\mu)$ of the log normalizing function $\Phi(\theta)$

$$(6) \quad (\text{cov}_\theta\{\phi(X)\})^{-1} = \nabla^2\Phi^*(\mu)$$

where the conjugate dual Φ^* of a function Φ is defined as

$$(7) \quad \Phi^*(\mu) := \sup_{\theta \in \Omega} \{\langle \mu, \theta \rangle - \Phi(\theta)\},$$

where $\mu \in \mathbb{R}^d$ is a fixed vector of mean parameters of the same dimension as θ , $\langle \mu, \theta \rangle$ is the Euclidean inner product $\sum_{i=1}^m \mu_i \theta_i$ and $\mu \in \mathcal{M}$, where

$\mathcal{M} := \{\mu \in \mathbb{R}^d \mid \exists P \text{ s.t. } \mathbb{E}_P[\phi_\alpha(X)] = \mu\}$. We define the mean parameter μ_α associated with a sufficient statistic ϕ_α as its expectation with respect to $P_\theta(x)$

$$(8) \quad \mu_\alpha := \mathbb{E}_P[\phi_\alpha(X)].$$

Note that, by assumption, $\Phi(\theta)$ is finite, which implies that \mathcal{M} is bounded. For an excellent treatment on the conjugate dual of the log normalizing function we refer the reader to [Wainwright and Jordan \(2008\)](#).

[Loh and Wainwright \(2013\)](#) showed that equality (6) holds for regular, minimal exponential family distributions with a finite log-normalizing constant $\Phi(\theta)$, that is, $\theta \in \Omega = \{\theta : \Phi(\theta) < \infty\}$ and we refer the reader to their paper for the proof of this result. Using this result, for the first part of the proof it remains to be shown that these requirements hold for the class of mixed distributions as in (4). [Yang et al. \(2014\)](#) show that the class is regular and has a finite log-normalizing constant, and minimality is established by Lemma 1:

LEMMA 1. *The class of mixed distributions as in 4 is minimal.*

For the proof of Lemma 1 see Appendix A.1.

The second part of the proof is to show that the graph structure claimed in Theorem 1 holds in $\nabla^2\Phi^*(\mu)$. In general there is no straight-forward way to compute $\nabla^2\Phi^*(\mu)$, because we do not know how to compute derivatives of $\Phi(\theta)$ in (7), which depends on all cliques in the graph. However, because we require the graph to be triangulated, we can represent P_θ as a factorization of marginal distributions of cliques $C \in \mathcal{C}$ and separator sets $S \in \mathcal{S}$ ([Lauritzen, 1996](#); [Cowell et al., 2007](#))

$$(9) \quad P_\theta = \frac{\prod_{C \in \mathcal{C}} P_C(X_C)}{\prod_{S \in \mathcal{S}} P_S(X_S)}.$$

If we plug the mixed joint distribution form as in (4) into (9) it is explicit that each term in the factorization depends only on the corresponding clique or separator set

$$(10) \quad P_\theta = \frac{\prod_{C \in \mathcal{C}} \exp \left\{ \frac{\theta_C \prod_{j=1}^{|C|} \phi_{t_j}(X_j) + \sum_{r \in C} C_r(X_r)}{\log \int_{x \in \mathcal{X}^{|C|}} \theta_C \prod_{j=1}^{|C|} \phi_{t_j}(X_j) + \sum_{r \in C} C_r(X_r)} \right\}}{\prod_{S \in \mathcal{S}} \exp \left\{ \frac{\theta_S \prod_{j=1}^{|S|} \phi_{t_j}(X_j) + \sum_{r \in S} C_r(X_r)}{\log \int_{x \in \mathcal{X}^{|S|}} \theta_S \prod_{j=1}^{|S|} \phi_{t_j}(X_j) + \sum_{r \in S} C_r(X_r)} \right\}}.$$

We obtain $\Phi^*(\mu)$ of $\Phi(\theta)$ via its equality relation to the negative Shannon entropy $H(P_\theta)$ of the distribution P_θ

$$\begin{aligned}
 \Phi^*(\mu) &= -H(P_\theta(\mu)) = -\mathbb{E}[\log P_\theta] = \\
 (11) \quad & - \sum_{C \in \mathcal{C}} \left(\theta_C \mu_{j \in C} + \sum_{r \in C} C_r(\mu_r) - \Phi_C(\theta) \right) \\
 & - \sum_{S \in \mathcal{S}} \left(\theta_S \mu_{j \in S} + \sum_{r \in S} C_r(\mu_r) - \Phi_S(\theta) \right).
 \end{aligned}$$

Note, that the factorization (10) turns into a sum in the conjugate dual representation (11), which enables to verify the claims of Theorem 1 by taking partial derivatives:

consider two subsets $A, B \in \mathcal{C}$ that are not contained in the same maximal clique C . As $P_\theta(X)$ is Markov with respect to the triangulated graph G , all interaction parameters $\theta_{A, \dots, B}$ associated with both A and B are equal to zero. When differentiating expression (11) with respect to μ_A all terms that do not involve μ_A drop. Now, when taking the second derivative with respect to μ_B , all terms that do not involve μ_B drop. The only terms left are terms involving both μ_A and μ_B . However, these terms are products involving a $\theta_{A, \dots, B} = 0$ and therefore we obtain $\frac{\partial^2 \Phi^*(\mu)}{\partial \mu_A \partial \mu_B} = 0$. Together with the equality (6), this proves part (1) of the graph structure in Theorem 1.

Turning to part (2), if A, B are part of the same maximal clique we have $\theta_{A, \dots, B} \neq 0$ for at least one parameter $\theta_{A, \dots, B}$. Taking partial derivatives with respect to μ_A and μ_B preserves only terms involving both μ_A and μ_B . Assuming the θ 's are drawn from a continuous distribution, $\frac{\partial^2 \Phi^*(\mu)}{\partial \mu_A \partial \mu_B}$ is almost surely nonzero. Together with equality (6), this proves part (2) of Theorem 1.

2.4.2. Proof of Corollary 1. We take the conjugate dual $\Phi^*(\mu)$ of $\Phi(\theta)$ corresponding to a model of the form (4) including only terms for cliques $(\{s\} \cup \mathbf{pow}(\mathcal{A}(s; d)))$. As above, we take derivatives with respect to μ_s and μ_B and all terms drop that do not contain both μ_s and μ_B . Whenever $\{s\} \neq B \not\subseteq N(s)$, all terms including both μ_s and μ_B are equal to zero, because $P_\theta(X)$ is Markov with respect to G . Therefore, whenever $\{s\} \neq B \not\subseteq N(s)$, the partial derivative of $\Phi^*(\mu)$ with respect to μ_s and μ_B is equal to zero. With equality (6), this proves the claims of Corollary 1.

2.5. Nodewise regression algorithm for structure estimation of mixed graphical models. In this section we describe how to use neighborhood-regression

together with the result in Corollary 1 to estimate a mixed MRF from data. Recall that we obtain the whole graph $G = (V, E)$ by recovering the neighborhood $N(s)$ of all $s \in V$. Corollary 1 enables us to construct an inverse covariance matrix Γ that reflects the neighborhood $N(s)$. We can therefore estimate the neighborhood of s by estimating the row s in Γ . Recall that this row s is a scalar multiple of the regression vector of X_s upon $X_{\setminus s}$ (see e.g. [Lauritzen, 1996](#)).

This means that we construct a random vector $\Psi(\phi(X); \{s\} \cup \text{pow}(\mathcal{A}(s; d)))$ for each node s as in Corollary 1 and we then predict the node s Ψ_s by all other terms $\Psi_{\setminus s}$ in Ψ . In order to get a sparse parameter vector $\hat{\theta}$, we apply ℓ_1 -regularization, whose strength is controlled by the regularization parameter λ_n . The parameter vector $\hat{\theta}$ is estimated by maximizing the ℓ_1 -penalized log likelihood. In case node s is associated with a Gaussian random variable, this is equivalent to solving the ℓ_1 -penalized least-squares problem

$$(12) \quad \hat{\theta} = \arg \min_{\|\theta\|_1 \leq b_0 \sqrt{k}} \{ \|\Psi_s - \Psi_{\setminus s} \theta\|_2 + \lambda_n \|\theta\|_1 \},$$

where $b_0 > \|\tilde{\theta}\|_1$ is a constant and k is the sparseness of the population parameter vector $\tilde{\theta}$. For non-Gaussian variables a (link) function is used to define a linear relation between the dependent variable Ψ_s and its predictors $\Psi_{\setminus s}$ (for details see [Friedman, Hastie and Tibshirani, 2010](#)).

Solving the lasso-problem in (12) for all $s \in V$ yields two estimates $\hat{\theta}_{i,j}$ and $\hat{\theta}_{j,i}$ for each edge. We combine these estimates with an AND-rule (both parameters have to be nonzero, otherwise the edge is absent; the value is the average) or an OR-rule (at least one parameter has to be nonzero; the parameter value is the average over the non-zero parameters).

Further, we assign the parameters λ_n, τ_n to the scaling

$$(13) \quad \lambda_n \asymp \sqrt{d} \|\tilde{\theta}\|_2 \sqrt{\frac{\log p}{n}}, \quad \tau_n \asymp \sqrt{d} \|\tilde{\theta}\|_2 \sqrt{\frac{\log p}{n}},$$

where $\|\tilde{\theta}\|_2$ is the ℓ_2 -norm of the population parameter vector $\tilde{\theta}$ and d is the degree of the true graph. By \asymp and \asymp we mean asymptotically larger and asymptotically over, respectively. Note the inequalities (13) contain the population parameter vector $\tilde{\theta}$. In order to compute τ_n , we use the estimated parameter vector $\hat{\theta}$ and we select λ_n using 10-fold cross-validation. We thereby obtain the following nodewise algorithm for structure estimation:

ALGORITHM 1. (*Nodewise regression method*)

1. Construct the vector Ψ according to Corollary 1.
2. Select a parameter λ_n using 10-fold cross-validation
3. Solve the lasso regression problem (12) with parameter λ_n , and denote the solution by $\hat{\theta}$.
4. Threshold the entries of $\hat{\theta}$ at level τ_n
5. Combine the neighborhood estimates with the AND- or OR-rule and define the estimated neighborhood set $\widehat{N}(s)$ as the support of the thresholded vector

As we do not know the degree of the true graph d , we have to make an assumption about d . The straightforward choice would be $d = |V| - 1$ which puts no constraints on the true graph. However, note that the computational complexity of Algorithm 1 is $\mathcal{O}(2^d \log p)$, which means choosing $d = |V| - 1$ renders the algorithm unfeasible for all but small graphs.

In order to make an informed assumption about d , note that the choice of d has two consequences in Algorithm 1: first, it influences the threshold τ_n in (13), which reflects the sparsity assumption necessary for asymptotic consistency. Second, d determines the size of candidate neighborhoods added to the covariance matrix (see Corollary 1). If we now only consider satisfying the requirements of Corollary 1, we can choose d as the size of the largest clique in graph G that does not contain s . For example, in the case that the true model is a pairwise model, we would choose $d = 2$ instead of the degree d of G , which depends on the size and density of graph G . The cost of the assumption about d is, with either interpretation, that we miss edges belonging to cliques with size $|C| > d$.

For the sample complexity $n \gtrsim d^3 \log p$, where d is the maximal degree in the true graph G , Loh and Wainwright (2013) prove asymptotic consistency for Algorithm 1. For details we refer the reader to their paper.

Because asymptotic considerations provide little information on how well a method performs in practical situations, we use simulations to test the performance of our method in recovering the graph from data in settings that resemble typical situations in exploratory data analysis.

3. Simulations. We consider the following six graphical models: Binary-Gaussian, Binary-Poisson, Binary-Exponential, and Multinomial with $m = 2, 3$, or 4 categories. In the three mixed graphical models, half of the nodes are binary. We used the Potts model (see e.g. Wainwright and Jordan, 2008) to specify interactions between categorical variables: if we have m categories, X_s, X_t are two arbitrary nodes and $j, k \in m$, then $\theta_{st,jj} = 1$ and $\theta_{st,jk} = 0$,

where $j \neq k$. Interactions between binary variables $X_s \in \{0, 1\}$ and continuous variables X_t are specified as follows: if $X_s = 1$ then $\theta_{st} = 1$ and if $X_s = 0$ then $\theta_{st} = 0$. All graphs were random graphs (Erdos and Rnyi, 1959) with $p = 16$ nodes and we varied sparsity by varying edge-probabilities $P_{edge} \in \{.1, .2, .3\}$. We also varied the ratio between observations and variables $\frac{n}{p} \in \exp\{0, 1, 2, 3, 4, 5\} \approx \{1, 3, 7, 20, 55, 148\}$ and the size of augmented cliques (interactions) $d \in \{1, 2, 3\}$, where $d = 1$ means that we use the original covariance matrix. Noise was added by multiplying each edge-weight by a draw from a uniform distribution $\mathcal{U}(.3, 1)$. All thresholds were set to zero, only for Poisson and Exponential nodes we set the threshold to $-.1$ in order to avoid $\lambda = 0$ in case $N(s) = \{\emptyset\}$. All edge-weights were additionally multiplied by -1 , which led to less extreme category-potentials and therefore more variance in the categorical variables.

In the Binary-Gaussian case, the variance of the Gaussian univariate distributions was set to 1 and the whole random graph including noise was resampled until the weighted adjacency matrix of the Gaussian sub-graph was positive definite. For categorical variables, we required that each category is present in the data and that the category with the lowest frequency is present in more than 10% of the cases. The reason for the second condition is that we need non-zero variance in every possible response vector in 10-fold cross-validation. For the same reason, we required for Poisson variables that the category with the highest frequency is present in less than 90% of the cases. In order to meet this requirement, we take the columns that do not meet the requirements and replace random row-entries by samples from categories with equal probabilities (categorical case) or by a draw from a Poisson distribution with $\lambda = \frac{1}{2}$ (Poisson case) until the requirement is met.

We combined parameters specifying the interaction between two categorical variables with the OR-rule and used the AND-rule to combine estimates between regressions.

Sensitivity and precision of our method for different combinations of sparsity (P_{edge}) and n/p -ratio for graphs consisting of categorical random variables with $m = \{2, 3, 4\}$ categories is shown in Figure 2. In all conditions with $d = 1$, sensitivity quickly converges to 1 with increasing observations. However, in this setting precision does not converge to 1. This is consistent with our theory, as we require in Corollary 1 that all candidate neighborhoods with size up to larger or equal to the largest clique in the true graph are augmented. As we sampled data from a pairwise model, the largest clique in the true graph has size $d = 2$ and we violate this requirement. For $k = 2, 3$ we satisfy the requirements in Corollary 1 and in these settings Algorithm 1 converges in precision. The general performance drops when the number

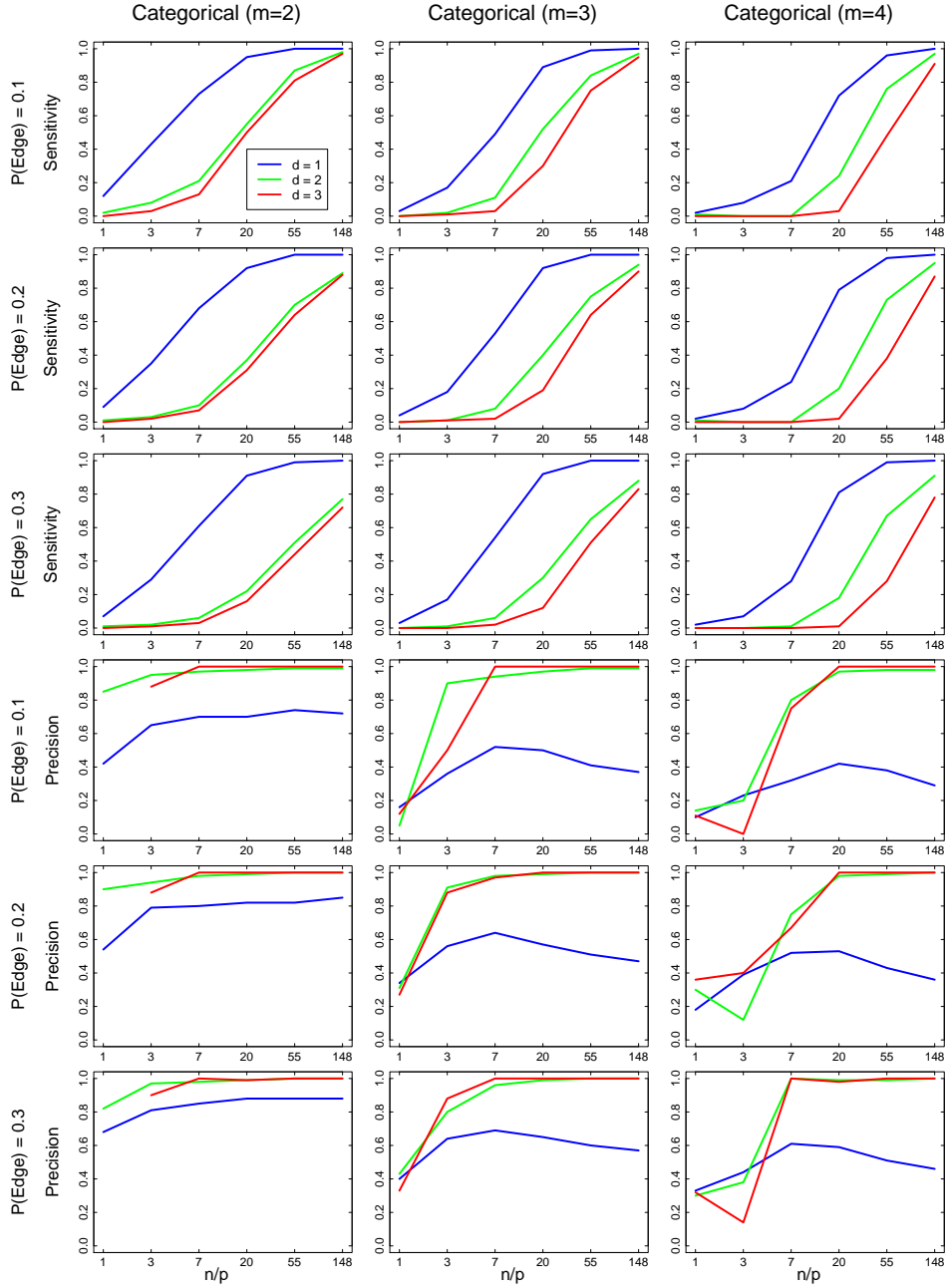


FIG 2. Simulation results: The first three rows show sensitivity, the last three rows precision of our method. We augmented cliques of size $d = \{1, 2, 3\}$ for different combinations of sparsity and rescaled sample size n/p for graphs consisting of categorical variables with $m = 2, 3$ or 4 . Missing data points in precision are due to the fact there were no edges estimated in any of the 100 iterations.

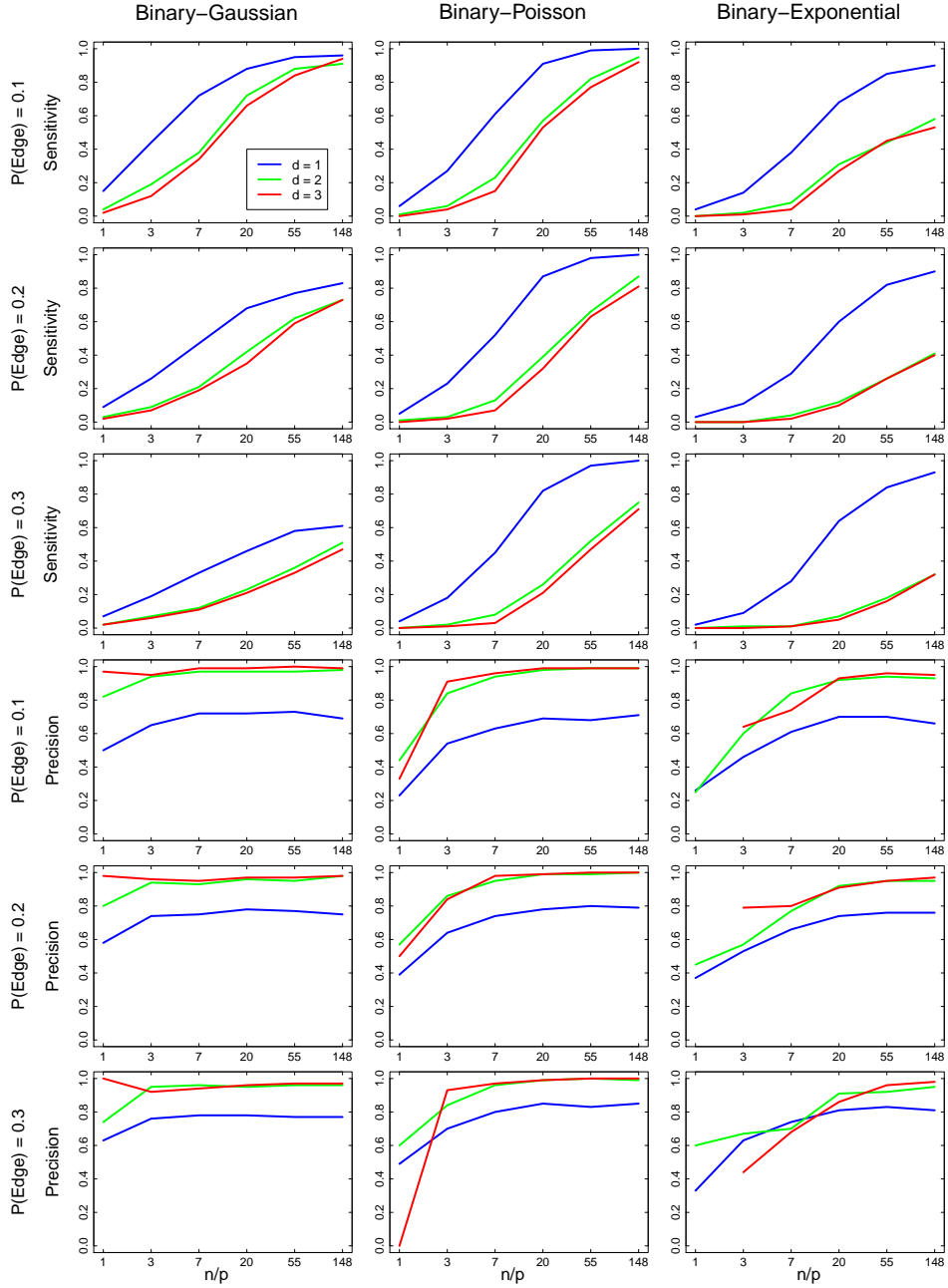


FIG 3. *Simulation results: The first three rows show sensitivity, the last three rows precision of our method. We augmented cliques of size $d = \{1, 2, 3\}$ for different combinations of sparsity and rescaled sample size n/p for Binary-Gaussian, Binary-Poisson and Binary-Exponential graphs. Missing data points in precision are due to the fact there were no edges estimated in any of the 100 iterations.*

of categories m increases, because more parameters have to be estimated. We also observe that general performance drops with lower sparsity, however the effects are small. Note that for the binary case $m = 2$ we obtain qualitatively similar results compared to [Loh and Wainwright \(2013\)](#) and others who used a similar estimation method ([Ravikumar, Wainwright and Lafferty, 2010](#); [van Borkulo et al., 2014](#)).

Note that we observe a peculiar pattern in precision in the categorical MRFs with $m > 2$ categories. Precision increases initially with observations, but then *decreases* again. This is a consequence of the OR-rule, which we use to combine several parameters describing the interaction between two categorical variables into one edge-parameter. As only one of the $(m - 1)^2$ parameters has to be nonzero to render the corresponding edge nonzero, the algorithm becomes more liberal when the number of categories is high. While these spurious parameters are put to zero by the τ_n -threshold (13) for small n , the threshold decreases with increasing n such that these spurious parameters are not set to zero anymore.

Figure 3 shows sensitivity and precision of our method for the same combinations of sparsity and n/p -ratio as above for Binary-Gaussian, Binary-Poisson and Binary-Exponential. Similarly to the categorical case, we see in conditions with $d = 1$ that sensitivity seems to converge to 1 with increasing n in all conditions but the Binary-Gaussian cases with sparsity $> .1$. In conditions with $d = 2, 3$ sensitivity seems to converge to 1 in the Binary-Gaussian case with sparsity $P_{edge} = .1$. In all other conditions with $d = 2, 3$ sensitivity increases with growing n . Precision converges most quickly for the Binary-Gaussian case, followed by the Binary-Poisson and Binary-Exponential case. Sparsity seems to have no impact on the precision in any of the three mixed graphs. Similarly to the categorical case, precision converges to 1 despite the fact that the theoretical sample complexity is not satisfied. As above, missing data points in precision is due to the fact that the algorithm often estimates zero edges.

4. Application to ASD Data. In this section we return to the exploratory data analysis problem introduced at the beginning of the paper. We estimate the MRF underlying a dataset consisting of responses to a questionnaire of 3521 individuals from the Netherlands diagnosed with Autism Spectrum Disorder (ASD). The variables cover demographics, social environment, diagnostic measurements and aspects of well-being (for details see [Begeer, Wierda and Venderbosch, 2013](#)).

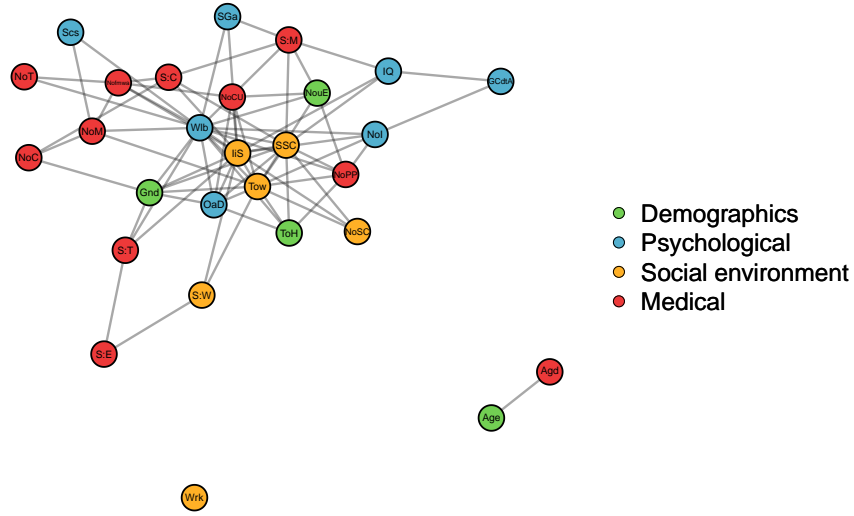
We assumed that the maximal clique size in the true graph is two and we therefore augment second-order interactions ($k = 2$) to the covariance ma-

trix and to combine parameters across regressions with the AND-rule. Different to the algorithm used for the simulations, we select the penalization-parameter λ using the EBIC, because the cross-validation requirement that each category is present in each fold was not satisfied in this dataset. For the layout in Figure 4, we used the force-directed algorithm of [Fruchterman and Reingold \(1991\)](#), that places nodes with many edges close to the center of the layout.

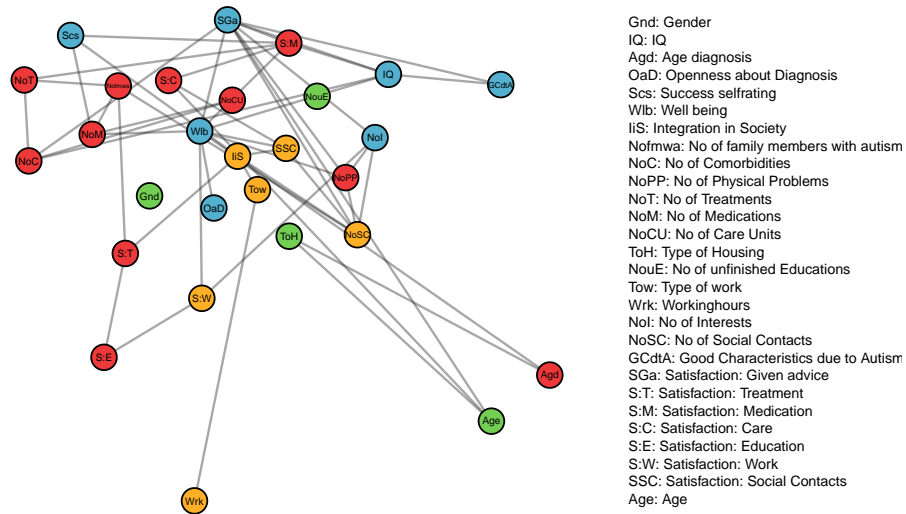
In Figure 4 (a), we see that the different aspects 'Demographics', 'Psychological', 'Social environment', and 'Medical' are strongly interrelated, which highlights the need for an integrated analysis. On the level of single variables, we can visually identify the importance of single variables, for example we see that 'Well-being' is central in the graph and has unique associations with many other variables. We can make the analysis of single variables more explicit by computing centrality measures for each node. Centrality measures quantify the importance of a node in a network, where the exact interpretation of importance depends on the specific centrality measure. In Figure 5 (a) we report the standardized centrality measures betweenness, closeness and degree (for details see [Opsahl, Agneessens and Skvoretz, 2010](#)). We see, for example, that 'Integration in Society' scores relatively high on closeness, degree and betweenness. This means 'Integration in Society' is relatively close to other nodes (closeness), has many connections (degree), and is often on the shortest path between any two nodes (betweenness).

We estimated the graph in Figure 4 (a) by modeling categorical variables as categorical variables, real-valued variables as Gaussian and count-valued variables as Poisson in a mixed graphical model. But does this actually give us a different graph compared with when we treat all variables as Gaussians? Figure 4 (b) shows the graph resulting from using the same estimation method as in (a), but treating all variables as Gaussians. We see that the resulting graph has less edges (density = .13 vs. .19). However, method (b) is not simply more liberal. We see also edges that are present in (a) but not in (b) such as the edge between 'Type of housing' and 'Openness about Diagnosis'. Also by comparing the centrality plots in 5 (a) and (b) we see considerable differences. For example, 'Satisfaction with Social Contacts' is one of the most central nodes in the mixed graphical model, while in the Gaussian graphical model its centrality is average or below average. These substantial differences between the two graphs highlight the importance of modeling variables on their proper domain.

5. Discussion. In the present paper, we extended the generalized covariance method of [Loh and Wainwright \(2013\)](#) to the class of mixed distri-

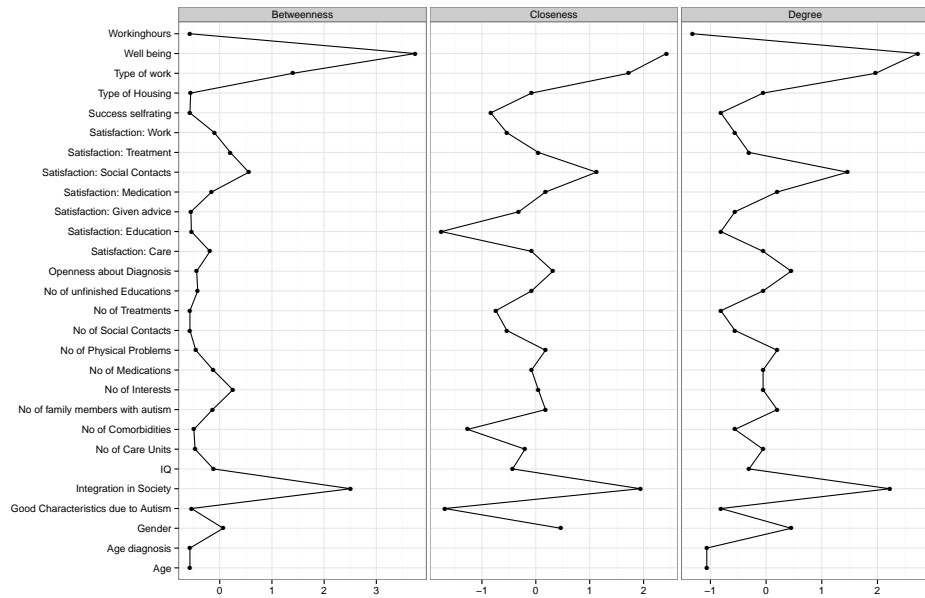


(a) Mixed Graphical Model

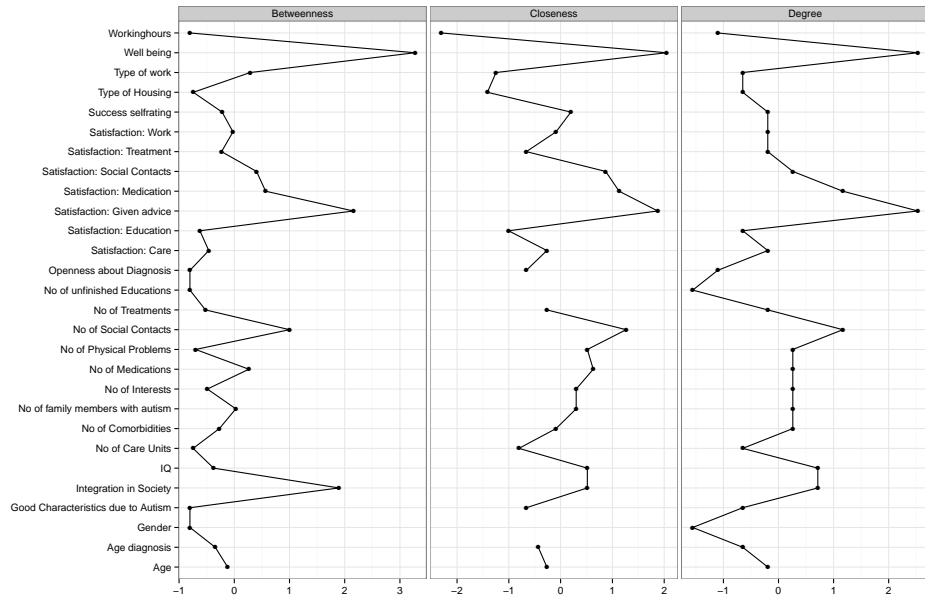


(b) Gaussian Graphical Model

FIG 4. Comparing two graphical models: (a) all variables are modeled on their proper domain (b) all variables are modeled as Gaussians.



(a) Mixed Graphical Model



(b) Gaussian Graphical Model

FIG 5. Comparing the centrality measures betweenness, closeness and degree of the mixed graphical model (a) and the Gaussian graphical model (b) in Figure 4.

butions introduced by Yang et al. (2014). We thereby provide an estimation method for the underlying graph structure of joint distributions over any combination of univariate members of the exponential family.

The peculiar simulation results of increasing and then decreasing precision with growing n in conditions with $m = 2, 3$ and $d = 1$ show that it is important to make sensible choices about how to combine the set of parameters involved in an interaction including a categorical variable into one dependence parameter. One solution to this problem might be using the group lasso (Yuan and Lin, 2006; Jacob, Obozinski and Vert, 2009), which induces sparse estimation on the group level.

As mentioned in Section 2.2, a limitation of the class of mixed distributions that in case they include two univariate distributions with infinite domain, they are not normalizable if *neither* both distributions are infinite only from one side *nor* the base measures are bounded with respect to the moments of the random variables. While this is an important theoretical problem that has to be addressed in future research, it does not limit the applicability of our method in most situations. Problems would only occur in the presence of extremely large values, which are rarely found in real-world data as most measurement scales are (naturally) bounded.

As indicated in Section 3, we added additional noise after sampling to be able to perform 10-fold cross-validation to select an appropriate λ_n parameter for the ℓ_1 -penalty in Algorithm 1. The proportion of data points replaced by noise until the requirement for cross-validation was met is depicted in Figure 6 in Appendix B. As can be seen in the figure, the proportion of added noise is considerable for the Binary-Gaussian and Binary-Exponential for small $\frac{n}{p}$ and in categorical graphs with $m = 2, 3$ and $\frac{n}{p} = 1$. Therefore, in cases with a nonzero proportion of additionally added noise, the above simulation-results can be interpreted as *conservative* estimates of how well the method performs in practice.

In order to make our method available to researchers, we implemented our method as the R-package **mgm**, that is available from the Comprehensive R Archive Network (CRAN) at <http://CRAN.r-project.org/>.

Statistical analyses of multivariate data based on Markov random fields are becoming increasingly popular in many areas of science. Despite the fact that most datasets involve different types of variables, up until now, there was no principled method available to estimate the Markov random field underlying a joint distribution over different types of variables. We closed this methodological gap by providing a well-performing and easy to interpret method to estimate the underlying Markov random field of a joint distribution consisting of any combination of univariate exponential family

members, which includes commonly used distributions such as Gaussian, Bernoulli, multinomial, Poisson, exponential, gamma, chi-squared and beta. We provided simulation results illustrating the performance of the method in realistic situations, illustrated our method with a network of different live aspects of individuals diagnosed with ASD and presented an implementation of our method as an R-package.

References.

- ANDERSON, D. K., LIANG, J. W. and LORD, C. (2014). Predicting young adult outcome among more and less cognitively able individuals with autism spectrum disorders. *Journal of Child Psychology and Psychiatry* **55** 485–494.
- BANERJEE, O., EL GHAOU, L. and D’ASPREMONT, A. (2008). Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. *The Journal of Machine Learning Research* **9** 485–516.
- BEGEER, S., WIERDA, M. and VENDERBOSCH, S. (2013). Allemaal Autisme, Allemaal Anders. Rapport NVA enquete 2013 [All Autism, All Different. Dutch Autism Society Survey 2013]. *Bilthoven: NVA*. 83.
- BORSBOOM, D. and CRAMER, A. O. J. (2013). Network Analysis: An Integrative Approach to the Structure of Psychopathology. *Annual Review of Clinical Psychology* **9** 91–121.
- CHENG, J., LEVINA, E. and ZHU, J. (2013). High-dimensional Mixed Graphical Models. *arXiv preprint arXiv:1304.2810*.
- COWELL, R. G., DAWID, P., LAURITZEN, S. L. and SPIEGELHALTER, D. J. (2007). *Probabilistic Networks and Expert Systems: Exact Computational Methods for Bayesian Networks*. Springer Science & Business Media.
- DOBRA, A. and LENKOSKI, A. (2011). Copula Gaussian graphical models and their application to modeling functional disability data. *The Annals of Applied Statistics* **5** 969–993.
- ERDOES, P. and RNYI, A. (1959). On random graphs, I. *Publicationes Mathematicae (Debrecen)* **6** 290–297.
- EVGENIOU, A. and PONTIL, M. (2007). Multi-task feature learning. *Advances in neural information processing systems* **19** 41.
- FOYGEL, R. and DRTON, M. (2014). High-dimensional Ising model selection with Bayesian information criteria. *arXiv preprint arXiv:1403.3374*.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9** 432–441.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software* **33** 1.
- FRUCHTERMAN, T. M. J. and REINGOLD, E. M. (1991). Graph drawing by force-directed placement. *Software: Practice and Experience* **21** 1129–1164.
- HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2009). *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York, New York, NY.
- HSU, C.-C., CHEN, C.-L. and SU, Y.-W. (2007). Hierarchical clustering of mixed data based on distance hierarchy. *Information Sciences* **177** 4474–4492.
- JACOB, L., OBOZINSKI, G. and VERT, J.-P. (2009). Group Lasso with Overlap and Graph Lasso. In *Proceedings of the 26th Annual International Conference on Machine Learning*. ICML ’09 433–440. ACM, New York, NY, USA.
- KOLLER, D. and FRIEDMAN, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT Press, Cambridge, MA.

- LAFFERTY, J., LIU, H. and WASSERMAN, L. (2012). Sparse Nonparametric Graphical Models. *Statistical Science* **27** 519–537.
- LAURITZEN, S. L. (1996). *Graphical models. Oxford statistical science series* **17**. Clarendon Press ; Oxford University Press, Oxford, New York.
- LEE, J. D. and HASTIE, T. J. (2012). Learning Mixed Graphical Models. *arXiv:1205.5012 [cs, math, stat]*. arXiv: 1205.5012.
- LIU, H., LAFFERTY, J. and WASSERMAN, L. (2009). The nonparanormal: Semiparametric estimation of high dimensional undirected graphs. *The Journal of Machine Learning Research* **10** 2295–2328.
- LIU, H., HAN, F., YUAN, M., LAFFERTY, J., WASSERMAN, L. and OTHERS (2012). High-dimensional semiparametric Gaussian copula graphical models. *The Annals of Statistics* **40** 2293–2326.
- LOH, P.-L. and WAINWRIGHT, M. J. (2013). Structure estimation for discrete graphical models: Generalized covariance matrices and their inverses. *The Annals of Statistics* **41** 3022–3049.
- MAGIATI, I., TAY, X. W. and HOWLIN, P. (2014). Cognitive, language, social and behavioural outcomes in adults with autism spectrum disorders: A systematic review of longitudinal follow-up studies in adulthood. *Clinical Psychology Review* **34** 73–86.
- MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the Lasso. *The Annals of Statistics* **34** 1436–1462.
- OPSAHL, T., AGNEESSENS, F. and SKVORETZ, J. (2010). Node centrality in weighted networks: Generalizing degree and shortest paths. *Social Networks* **32** 245–251.
- RAVIKUMAR, P., WAINWRIGHT, M. J. and LAFFERTY, J. D. (2010). High-dimensional Ising model selection using l_1 -regularized logistic regression. *The Annals of Statistics* **38** 1287–1319.
- RUE, H., MARTINO, S. and CHOPIN, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **71** 319–392.
- SAMMEL, M. D., RYAN, L. M. and LEGLER, J. M. (1997). Latent Variable Models for Mixed Discrete and Continuous Outcomes. *Journal of the Royal Statistical Society: Series B (Methodological)* **59** 667–678.
- VAN BORKULO, C. D., BORSBOOM, D., EPSKAMP, S., BLANKEN, T. F., BOSCHLOO, L., SCHOEVERS, R. A. and WALDORP, L. J. (2014). A new method for constructing networks from binary data. *Scientific Reports* **4**.
- WAINWRIGHT, M. J. and JORDAN, M. I. (2008). Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends in Machine Learning* **1** 1–305.
- XUE, L. and ZOU, H. (2012). Regularized rank-based estimation of high-dimensional nonparanormal graphical models. *The Annals of Statistics* **40** 2541–2571.
- YANG, X., KIM, S. and XING, E. P. (2009). Heterogeneous multitask learning with joint sparsity constraints. In *Advances in Neural Information Processing Systems 22* (Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams and A. Culotta, eds.) 2151–2159. Curran Associates, Inc.
- YANG, E., BAKER, Y., RAVIKUMAR, P., ALLEN, G. and LIU, Z. (2014). Mixed Graphical Models via Exponential Families. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics* 1042–1050.
- YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68** 49–67.
- YUAN, M. and LIN, Y. (2007). Model selection and estimation in the Gaussian graphical model. *Biometrika* **94** 19–35.

APPENDIX A: PROOFS OF SUPPORTING LEMMAS

A.1. Proof of Lemma 1. The log-normalizing constant of the mixed joint distribution in (4) can be written as

$$\Phi(\theta) := \log \int_{x \in \mathcal{X}^p} \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C \phi(x_C) \right\} \nu(dx),$$

where ν is the count-measure, the Lebesgue-measure, or a combination of both.

To establish minimality, suppose $\sum_C a_C \phi(x_C) = b$ almost surely, where the coefficients a_C are real-valued and b is some constant. Plugging in x such that $x_s = 0$ for all $s \in V$ and using the fact that all states (discrete random variables) and intervals (continuous random variables) have positive probability, we see that $b = 0$. Now assume that not all a_C are equal to 0. Let C' be a set of cliques such that $a_{C'} \neq 0$ and $|C'|$ is minimal. Plugging in x such that $x_{C'} \neq 0$ and $x_{C \setminus C'} = 0$, we have

$$\sum_C a_C \phi(x_C) = a_{C'}$$

by the minimality of $|C'|$. This contradicts the fact $a_{C'} \neq 0$. Hence, we conclude that the sufficient statistics $\phi(x_C)$ are indeed linearly independent, implying that the class of mixed distributions as in (4) is minimal.

A.2. Corollary 2. Let $\mathbf{pow}(\mathcal{A}) = \bigcup_{C \in \mathcal{A}} \mathbf{pow}(C)$ be the union of all $2^{|C|} - 1$ nonempty subsets of all cliques C .

COROLLARY 2. *Take any triangulation \tilde{G} of the graph G and let \mathcal{A} be the set of separator sets in \tilde{G} . Then the inverse Γ of the covariance matrix $\text{cov}(\Psi(\phi(X); ; V \cup \mathbf{pow}(\mathcal{A})))$ has the property that $\tilde{\Gamma}(\{s\}, \{t\}) = 0$ whenever $(s, t) \notin \tilde{E}$.*

A consequence of Corollary A.2 is that for all graphs with singleton separator sets (e.g. trees), the original covariance matrix $\text{cov}(\Psi(\phi(X), V))$ is graph structured.

Proof. The proof of Corollary A.2 follows directly from the proof of Theorem 1. We take the conjugate dual $\Phi^*(\mu)$ as in (11) but of the log normalization function $\Phi(\theta)$ corresponding to a model of the form (4) including only terms for cliques $(V \cup \mathbf{pow}(\mathcal{A}))$. Analogously to the proof of Theorem 1 we take derivatives with respect to μ_s and μ_t and thereby all terms drop that do not contain both μ_s and μ_t . Whenever $(s, t) \notin E$, all terms including both μ_s and μ_t are equal to zero, because $P_\theta(X)$ is Markov with respect to G .

Therefore, whenever $(s, t) \notin E$, the partial derivative of $\Phi^*(\mu)$ with respect to μ_s and μ_t is equal to zero. With equality (6), this proves Corollary A.2.

APPENDIX B: ADDITIONAL NOISE PROPORTION IN SIMULATION

Figure 6 depicts the proportion of noise added after sampling the data until the requirements for 10-fold cross-validation were met. A proportion of 1 means that all data points were replaced by noise.

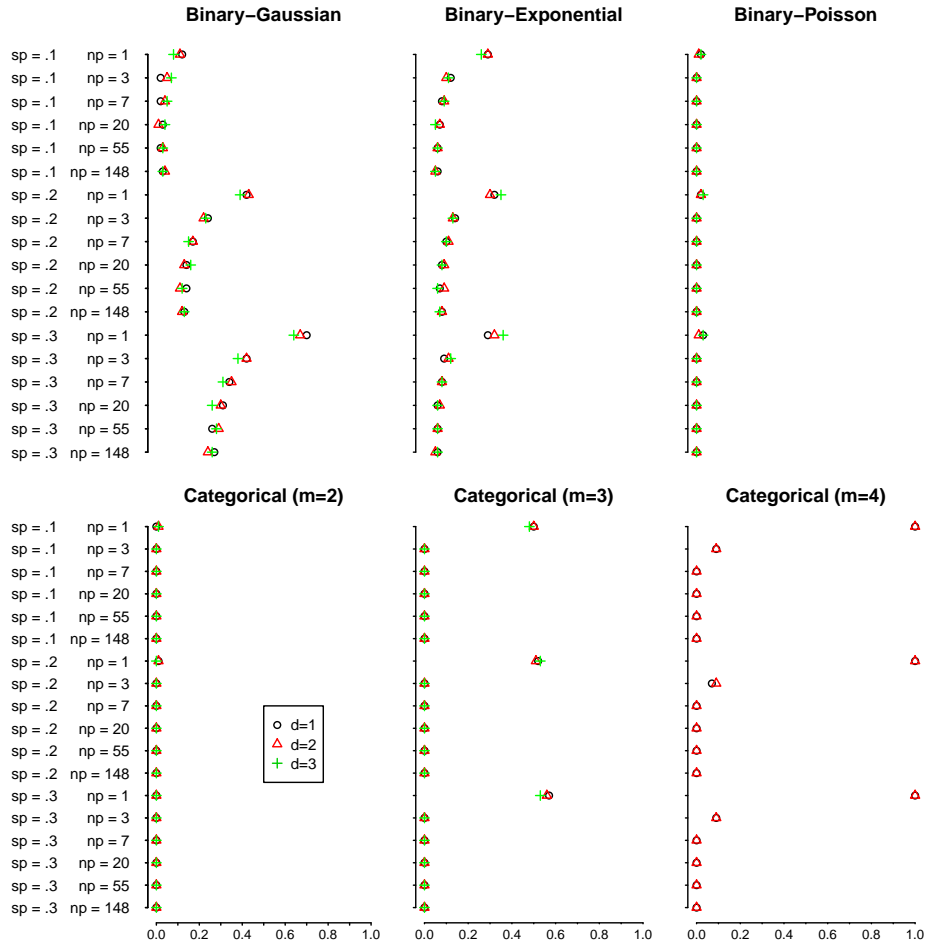


FIG 6. Proportion of additional noise added to the data until the requirements for 10-fold cross-validation were met.

E-MAIL: jonashasbeck@gmail.com

E-MAIL: waldorp@uva.nl