

Replicates in high dimensions, with applications to latent variable graphical models

BY KEAN MING TAN

*Department of Operations Research and Financial Engineering, Princeton University,
Princeton, New Jersey 08544, U.S.A.*

kmtan@princeton.edu

5

YANG NING

Department of Statistical Science, Cornell University, Ithaca, New York 14853, U.S.A.

yn265@cornell.edu

DANIELA M. WITTEN

Department of Statistics, University of Washington, Seattle, Washington 98195, U.S.A.

dwitten@uw.edu

10

AND HAN LIU

*Department of Operations Research and Financial Engineering, Princeton University,
Princeton, New Jersey 08544, U.S.A.*

hanliu@princeton.edu

15

SUMMARY

In classical statistics, much thought has been put into experimental design and data collection. In the high-dimensional setting, however, experimental design has been less of a focus. In this paper, we stress the importance of collecting multiple replicates for each subject in this setting. We consider learning the structure of a graphical model with latent variables, under the assumption that these variables take a constant value across replicates within each subject. By collecting multiple replicates for each subject, we are able to estimate the conditional dependence relationships among the observed variables given the latent variables. To test the null hypothesis of conditional independence between two observed variables, we propose a pairwise decorrelated score test. Theoretical guarantees are established for parameter estimation and for this test. We show that our proposal is able to estimate latent variable graphical models more accurately than some existing proposals, and apply the proposed method to a brain imaging dataset.

20

25

Some key words: Experimental design; Nuisance parameter; Pairwise decorrelated score test; Semiparametric exponential family graphical model.

30

1. INTRODUCTION

Experimental design and data collection have been the subjects of extensive research (Box et al., 2005; Montgomery, 2008). For instance, randomised clinical trials are conducted to determine the treatment effect of a new drug, and sample size calculations are performed to determine the smallest number of patients needed to give sufficient power to detect the treatment effect. In

35

contrast, in the high-dimensional setting, statisticians are usually not involved in experimental design and data collection. Given a cost constraint, investigators often try to obtain the largest possible number of subjects; that is, replicates are typically not collected for each subject. In this paper, we show that collecting replicates aids when learning an undirected graphical model with latent variables.

In an undirected graphical model, each node represents a random variable, and an edge connecting a pair of nodes indicates that the two variables are conditionally dependent, given all the other variables. The Gaussian graphical model has been studied extensively (Meinshausen & Bühlmann, 2006; Yuan & Lin, 2007; Friedman et al., 2008; Rothman et al., 2008; Peng et al., 2009; Ravikumar et al., 2011; Cai et al., 2011; Sun & Zhang, 2013). Other authors have considered extensions to the case in which each node-conditional distribution belongs to a univariate exponential family (Ravikumar et al., 2010; Yang et al., 2015; Lee & Hastie, 2015; Chen et al., 2015). Others have considered estimating conditional dependence relationships using semiparametric or nonparametric approaches (Liu et al., 2009, 2012; Fellinghauer et al., 2013; Voorman et al., 2014).

However, in many scientific studies, we observe only a subset of the relevant variables. For instance, in the context of a gene expression study, certain patients may have undiagnosed disease or some unknown risk factors. If the heterogeneity among patients is ignored, then the estimated conditional relationships among the genes may be distorted. This is made apparent in recent work on Gaussian graphical modelling in the presence of latent variables (Chandrasekaran et al., 2012), which showed that after marginalizing over the latent variables, the conditional independence graph corresponding to the observed variables may be dense.

In this paper, we propose an estimator and develop theory for the semiparametric exponential family graphical model with latent variables. This work builds upon an unpublished 2014 technical report by Yang et al. (arXiv:1412.8697), in which the semiparametric exponential family graphical model was introduced. We assume that these variables are constant across replicates within a given subject and that we have at least two replicates per subject. We exploit the replicates in order to construct a nuisance-free loss function that does not depend on the latent variables. In addition, we propose a pairwise decorrelated score test of the null hypothesis that two variables are conditionally independent, given all the other variables.

2. A MODEL FOR LATENT VARIABLE GRAPHICAL MODELS

2.1. Review of the semiparametric exponential family graphical model

We provide a brief review of the semiparametric exponential family graphical model proposed in Yang et al. (arXiv:1412.8697). Let X be a p -dimensional random vector and let $X_{-j} = (X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p)^T$. The p -dimensional random vector X follows the semiparametric exponential family graphical model if, for any node j , the conditional density of X_j given X_{-j} satisfies

$$p(x_j | x_{-j}) = \exp \{x_j \beta_{j,-j}^T x_{-j} + f_j(x_j) - A_j(\beta_j, f_j)\}, \quad (1)$$

where $\beta_{j,-j}$ encodes the conditional dependence relationships between the j th node and the other nodes, $f_j(\cdot)$ is an unknown function, and $A_j(\cdot)$ is the log-partition function. Because $f_j(\cdot)$ is unknown, obtaining the maximum likelihood estimator of (1) may be infeasible. To estimate $\beta_{j,-j}$, we can instead construct a loss function that does not depend on $f_j(\cdot)$.

Let X_i and x_i be the random variables and data corresponding to the i th subject, respectively. Let $x_{\cdot j} = (x_{1j}, \dots, x_{nj})^T$, and let $x_j^{(\cdot)}$ and $z_{\cdot j}$ be the order and rank statistics of $x_{\cdot j}$, respec-

tively. For instance, if $x_{.j} = (1, 5, 2)^\top$, then $x_j^{(\cdot)} = (1, 2, 5)^\top$ and $z_{.j} = (1, 3, 2)^\top$. Furthermore, let $x_{\cdot,-j}$ denote an $n \times (p-1)$ matrix obtained by stacking the vectors $x_{.k}$ for $k \neq j$. The joint conditional density of the j th variable given the others can be decomposed as

$$p(x_{.j} \mid x_{\cdot,-j}, \beta_{j,-j}, f_j) = p\{z_{.j} \mid x_{\cdot,-j}, x_j^{(\cdot)}, \beta_{j,-j}\} p\{x_j^{(\cdot)} \mid x_{\cdot,-j}, \beta_{j,-j}, f_j\},$$

the product of the conditional density of the rank statistics given the order statistics, and the density of the order statistics. The former does not depend on $f_j(\cdot)$: the key insight is that the rank statistics given the order statistics have no information about $f_j(\cdot)$. Rather than estimating $\beta_{j,-j}$ from the joint conditional density that involves the unknown function $f_j(\cdot)$, we can estimate $\beta_{j,-j}$ by maximising the conditional density of the rank statistics.

However, computing the conditional density of the rank statistics may be computationally prohibitive. Thus, we can consider the conditional density formed by a single pair of samples, and construct a nuisance-free likelihood function by multiplying the conditional densities of the $n(n-1)/2$ pairs of samples. This approach is also considered in Ning et al. (2016) in the context of semiparametric regression.

2.2. Semiparametric exponential family graphical models with latent variables

We generalise the semiparametric exponential family graphical model to accommodate latent variables. Let $X = (X_O^\top, X_H^\top)^\top$ be a $(p+h)$ -dimensional random vector, where $X_O \in \mathbb{R}^p$ and $X_H \in \mathbb{R}^h$ are the vectors of observed and latent random variables, respectively. We let $O = \{1, \dots, p\}$ and $H = \{p+1, \dots, p+h\}$ denote the index sets of the observed and latent random variables, respectively.

DEFINITION 1. A $(p+h)$ -dimensional random vector $X = (X_O^\top, X_H^\top)^\top$ follows a semiparametric exponential family graphical model with latent variables, if for any node j , the conditional density of X_j given X_{-j} satisfies

$$p(x_j \mid x_{-j}) = \exp \{x_j \beta_{j,-j}^\top x_{-j} + f_j(x_j) - A_j(\beta_j, f_j)\},$$

where $f_j(x_j)$ is some possibly unknown function and $A_j(\beta_j, f_j)$ is the log-partition function.

For any $j \in O$, we write $X_{O \setminus j} = (X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p)^\top \in \mathbb{R}^{p-1}$ and $X_{-j} = (X_{O \setminus j}^\top, X_H^\top)^\top \in \mathbb{R}^{p+h-1}$. Let $\beta_{j,O \setminus j} = (\beta_{j1}, \dots, \beta_{j,j-1}, \beta_{j,j+1}, \dots, \beta_{jp})^\top \in \mathbb{R}^{p-1}$, $\beta_{j,H} = (\beta_{j,p+1}, \dots, \beta_{j,p+h})^\top \in \mathbb{R}^h$, and $\beta_j = (\beta_{j,O \setminus j}^\top, \beta_{j,H}^\top)^\top$. The model introduced in Definition 1 can be rewritten as

$$p(x_j \mid x_{-j}) = \exp \left\{ x_j \beta_{j,O \setminus j}^\top x_{O \setminus j} + x_j \beta_{j,H}^\top x_H + f_j(x_j) - A_j(\beta_j, f_j) \right\}. \quad (2)$$

The parameters $\beta_{j,O \setminus j}$ and $\beta_{j,H}$ encode the conditional dependence relationships between the j th node and all the other observed and latent variables, respectively. In particular, $\beta_{jk} = 0$ if and only if the j th and k th nodes are conditionally independent, given all the other nodes.

In this paper, we assume that: $\beta_{jk} = \beta_{kj}$ and $\exp\{\sum_{j=1}^{p+h} \sum_{k \neq j} \beta_{jk} x_j x_k / 2 + \sum_{j=1}^{p+h} f_j(x_j)\}$ is integrable with respect to its measure. By an application of Proposition 1 in Chen et al. (2015), under these conditions, there exists a joint probability distribution for the model introduced in Definition 1 that takes the form

$$p(x) \propto \exp \left\{ \frac{1}{2} \sum_{j=1}^{p+h} \sum_{k \neq j} \beta_{jk} x_j x_k + \sum_{j=1}^{p+h} f_j(x_j) \right\}. \quad (3)$$

We provide two special cases of (2), and consider them in Section 4.

Example 1. The Gaussian graphical model with latent variables: let $X = (X_O^T, X_H^T)^T \sim N(0, \Sigma)$, where $\Sigma \in \mathbb{R}^{(p+h) \times (p+h)}$ and let $\Theta = \Sigma^{-1}$. For $j \in O$, the conditional density of X_j given all the other variables is

$$p(x_j | x_{-j}) = \left(\frac{\Theta_{jj}}{2\pi} \right)^{1/2} \exp \left\{ -x_j \Theta_{j,O \setminus j}^T x_{O \setminus j} - x_j \Theta_{j,H}^T x_H - \Theta_{jj} x_j^2 / 2 - \frac{(\Theta_{j,-j}^T x_{-j})^2}{2\Theta_{jj}} \right\}.$$

Comparing this to (2), we see that $\beta_{j,O \setminus j} = -\Theta_{j,O \setminus j}$, $\beta_{j,H} = -\Theta_{j,H}$, $f_j(x_j) = -\Theta_{jj} x_j^2 / 2$, and $A_j(\beta_j, f_j) = (\sum_{k \neq j} \Theta_{jk} x_k)^2 / (2\Theta_{jj}) + \log(2\pi/\Theta_{jj})/2$.

Example 2. The Ising model with latent variables: let $X_j \in \{0, 1\}$ with joint density $p(x) \propto \exp(\sum_{j < k} \Theta_{jk} x_j x_k)$. For $j \in O$, the conditional density of X_j given the other variables is

$$p(x_j | x_{-j}) = \exp \left[x_j \Theta_{j,O \setminus j}^T x_{O \setminus j} + x_j \Theta_{j,H}^T x_H - \log \{ 1 + \exp(\Theta_{j,-j}^T x_{-j}) \} \right].$$

Comparing this to (2), we see that $\beta_{j,O \setminus j} = \Theta_{j,O \setminus j}$, $\beta_{j,H} = \Theta_{j,H}$, $f_j(x_j) = 0$, and $A_j(\beta_j, f_j) = \log\{1 + \exp(\Theta_{j,-j}^T x_{-j})\}$.

2.3. From replicates to a nuisance-free loss function

Recall that we are interested in estimating the conditional dependence relationships among the observed variables given the latent variables, $\beta_{j,O \setminus j}$. Due to the presence of the possibly unknown function $f_j(x_j)$ and the latent variables x_H in (2), it is not possible to directly maximise (2) with respect to $\beta_{j,O \setminus j}$. By collecting multiple replicates per subject, and assuming that the latent variables are constant across replicates for a given subject, we construct a loss function that does not depend on the latent variables and the unknown function.

Let R_i be the number of replicates for the i th subject. To simplify bookkeeping, we assume that $R_1 = \dots = R_n = R$, though this assumption is not critical. Suppose that X_i^r , the random vector for the r th replicate for the i th subject is distributed as in (3), for $i = 1, \dots, n$ and $r = 1, \dots, R$. Throughout the paper, we assume that X_i^r and $X_{i'}^r$ are independent, while X_i^r and $X_i^{r'}$ may be dependent. Let $x_i^r = \{(x_{iO}^r)^T, (x_{iH}^r)^T\}^T$ be the data corresponding to the r th replicate of the i th subject. We make two assumptions on the replicates.

Assumption 1. The latent variables are constant across replicates, that is, $x_{iH}^r = x_{iH}^{r'} = x_{iH}$ for all $1 \leq r' \leq r \leq R$.

Assumption 2. Given the latent variables, the R replicates are mutually independent. That is, $p(x_{iO}^1, \dots, x_{iO}^R | x_{iH}) = \prod_{r=1}^R p(x_{iO}^r | x_{iH})$.

Assumptions 1 and 2 are plausible in many scientific settings. For instance, consider a gene expression study in which the expression levels of thousands of genes are measured for a number of subjects. Certain subjects may have unknown risk factors that might be associated with their gene expression levels. In this setting, the observed variables are the genes, and the latent variables represent unknown risk factors. Assumption 1 is satisfied if the disease status or the unknown risk factors do not change across time. Assumption 2 is likely to be satisfied, if the gene expression levels are measured in multiple independent clinical visits.

We now construct a nuisance-free loss function using an approach similar to the one outlined in Section 2.1, by exploiting the fact that R replicates are available for each subject. Under

Assumption 2, the joint conditional density for the i th subject for $j \in O$ is

$$p(x_{ij}^1, \dots, x_{ij}^R \mid x_{i,-j}^1, \dots, x_{i,-j}^R) = \prod_{r=1}^R p(x_{ij}^r \mid x_{i,-j}^r).$$

Estimating $\beta_{j,O \setminus j}$ by maximising the joint conditional density, which depends on both the unmeasured data x_{iH} and the possibly unknown function $f_j(\cdot)$, may not be feasible. 150

Let $x_{ij}^{(r,r')} = \{\min(x_{ij}^r, x_{ij}^{r'}), \max(x_{ij}^r, x_{ij}^{r'})\}$ be the order statistics of a pair of replicates for the i th subject. The joint conditional density for the pair of replicates is

$$\begin{aligned} p(X_{ij}^r = x_{ij}^r, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^r, x_{i,-j}^{r'}) \\ = p\{X_{ij}^r = x_{ij}^r, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^r, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\} p\{x_{ij}^{(r,r')} \mid x_{i,-j}^r, x_{i,-j}^{r'}\}. \end{aligned} \quad (4)$$

The following proposition shows that the conditional density $p\{X_{ij}^r = x_{ij}^r, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^r, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\}$ does not depend on $x_{iH}^r, x_{iH}^{r'}$, or on the unknown function $f_j(\cdot)$. 155

PROPOSITION 1. Under Assumptions 1 and 2, for $j \in O$,

$$p\{X_{ij}^r = x_{ij}^r, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^r, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\} = \left\{1 + R_{ij}^{rr'}(\beta_{j,O \setminus j})\right\}^{-1}, \quad (5)$$

where $R_{ij}^{rr'}(\beta_{j,O \setminus j}) = \exp\{-(x_{ij}^r - x_{ij}^{r'})\beta_{j,O \setminus j}^T(x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'})\}$.

In the absence of latent variables, similar results were established in the context of the semi-parametric generalised linear model (Equation (3.4) in Ning et al., 2016) and the semiparametric exponential family graphical model (Equation (3.1) in Yang et al. (arXiv:1412.8697)), both of which applied the approach outlined in Section 2.1. 160

Remark 1. When Assumption 1 is violated, the conditional density (5) takes the form

$$\left[1 + \exp\left\{-(x_{ij}^r - x_{ij}^{r'})\beta_{j,O \setminus j}^T(x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'}) - (x_{ij}^r - x_{ij}^{r'})\beta_{j,H}^T(x_{iH}^r - x_{iH}^{r'})\right\}\right]^{-1},$$

which has an additional term $(x_{ij}^r - x_{ij}^{r'})\beta_{j,H}^T(x_{iH}^r - x_{iH}^{r'})$ that depends on the latent variables. Provided that $|x_{iH}^r - x_{iH}^{r'}|$ is sufficiently small, this term is ignorable, and therefore it has negligible effect on the estimation of $\beta_{j,O \setminus j}$. 165

To obtain an estimate of $\beta_{j,O \setminus j}$, we ignore the term $p\{x_{ij}^{(r,r')} \mid x_{i,-j}^r, x_{i,-j}^{r'}\}$ in (4), and consider the product of joint conditional densities over all pairs of replicates across the n subjects,

$$\prod_{i=1}^n \prod_{1 \leq r < r' \leq R} p\{x_{ij}^r, x_{ij}^{r'} \mid x_{i,-j}^r, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\}.$$

This leads to a nuisance-free loss function that does not depend on the latent variables or on the unknown function, i.e.,

$$\ell_j(\beta_{j,O \setminus j}) = \frac{1}{n} \sum_{i=1}^n \left[\binom{R}{2}^{-1} \sum_{1 \leq r < r' \leq R} \log \left\{1 + R_{ij}^{rr'}(\beta_{j,O \setminus j})\right\} \right]. \quad (6)$$

From now onwards, we let $\beta_{j,O \setminus j}^*$ be the true parameter values in (2) that encode the underlying conditional dependence relationships between the j th variable and the observed variables. 170

Similarly, we let f_j^* be the underlying function in (2). The following proposition justifies the use of the loss function (6) for estimating $\beta_{j,O \setminus j}^*$.

PROPOSITION 2. For all $j \in O$, $E\{\nabla \ell_j(\beta_{j,O \setminus j}^*)\} = 0$ and $\beta_{j,O \setminus j}^*$ is a global minimizer of $E\{\ell_j(\beta_{j,O \setminus j})\}$, where $E(\cdot)$ is the expectation under the true parameters $(\beta_{j,O \setminus j}^*, f_j^*)$.

To encourage the estimated parameter to contain many zero elements, we solve

$$\underset{\beta_{j,O \setminus j} \in \mathbb{R}^{p-1}}{\text{minimize}} \quad \{\ell_j(\beta_{j,O \setminus j}) + \lambda \|\beta_{j,O \setminus j}\|_1\}, \quad (7)$$

where λ is a non-negative tuning parameter that controls the sparsity of the estimate $\hat{\beta}_{j,O \setminus j}$. The loss function (6) can be interpreted as a logistic loss function with $x_{ij}^r - x_{ij}^{r'}$ as the outcome and $x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'}$ as the covariates. We create a pseudo-binary outcome $\tilde{x}_{ij}^{rr'} = \text{sign}(x_{ij}^r - x_{ij}^{r'})$ and pseudo covariates $\tilde{x}_{i,O \setminus j}^{rr'} = (x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'})|x_{ij}^r - x_{ij}^{r'}|$. We can then solve (7) using the R package glmnet for fitting an ℓ_1 -penalised logistic regression to obtain an estimate of $\beta_{j,O \setminus j}$. When there are ties in the outcome, that is, $x_{ij}^r = x_{ij}^{r'}$, we ignore the pair of observations, since its contribution to the loss function (6) is free of the parameter of interest, $\beta_{j,O \setminus j}$.

2.4. Pairwise decorrelated score test

In this section, we consider testing a pre-specified component in $\beta_{j,O \setminus j}^*$ and $\beta_{k,O \setminus k}^*$, that is,

$$H_0 : \beta_{jk}^* = \beta_{kj}^* = 0 \quad \text{versus} \quad H_1 : \beta_{jk}^* = \beta_{kj}^* \neq 0, \quad (8)$$

for any $j, k \in O$, by treating the remaining parameters $\beta_{j,O \setminus \{j,k\}}^*$ and $\beta_{k,O \setminus \{j,k\}}^*$ as nuisance parameters. The classical score test is often used for this purpose in the low-dimensional setting. However, in the high-dimensional setting, the score test statistic is not asymptotically normal, because the number of nuisance parameters is large. We propose a pairwise decorrelated score test to test the null hypothesis given in (8). The test is constructed so that the effect of the nuisance parameters is asymptotically negligible. The decorrelated score test has been previously considered in Ning & Liu (2016), Ning et al. (2016), and Yang et al. (arXiv:1412.8697).

Let $\nabla \ell_j(\beta_{j,O \setminus j}) \in \mathbb{R}^{p-1}$ and $\nabla^2 \ell_j(\beta_{j,O \setminus j}) \in \mathbb{R}^{(p-1) \times (p-1)}$ be the gradient and the Hessian of the loss function $\ell_j(\beta_{j,O \setminus j})$ in (6), respectively. For $k \in O \setminus j$, we let

$$\nabla_k \ell_j(\beta_{j,O \setminus j}) = \frac{\partial \ell_j(\beta_{j,O \setminus j})}{\partial \beta_{jk}} \in \mathbb{R}, \quad \nabla_{-k} \ell_j(\beta_{j,O \setminus j}) = \frac{\partial \ell_j(\beta_{j,O \setminus j})}{\partial \beta_{j,O \setminus \{j,k\}}} \in \mathbb{R}^{p-2}.$$

Similarly, for $k \in O \setminus j$, we let

$$\nabla_{k,-k}^2 \ell_j(\beta_{j,O \setminus j}) = \frac{\partial^2 \ell_j(\beta_{j,O \setminus j})}{\partial \beta_{jk} \partial \beta_{j,O \setminus \{j,k\}}} \in \mathbb{R}^{p-2}, \quad \nabla_{-k,-k}^2 \ell_j(\beta_{j,O \setminus j}) = \frac{\partial^2 \ell_j(\beta_{j,O \setminus j})}{(\partial \beta_{j,O \setminus \{j,k\}})^2} \in \mathbb{R}^{(p-2) \times (p-2)}.$$

Define $H^j = E\{\nabla^2 \ell_j(\beta_{j,O \setminus j}^*)\} \in \mathbb{R}^{(p-1) \times (p-1)}$, and for $k \in O \setminus j$, let

$$H_{k,-k}^j = E\left\{\nabla_{k,-k}^2 \ell_j(\beta_{j,O \setminus j}^*)\right\} \in \mathbb{R}^{p-2}, \quad H_{-k,-k}^j = E\left\{\nabla_{-k,-k}^2 \ell_j(\beta_{j,O \setminus j}^*)\right\} \in \mathbb{R}^{(p-2) \times (p-2)}.$$

Let $(w_{jk}^*)^T = (H_{k,-k}^j)^T (H_{-k,-k}^j)^{-1} \in \mathbb{R}^{p-2}$, and let $\beta_{j \vee k} = (\beta_{jk}, \beta_{j,O \setminus \{j,k\}}^T, \beta_{k,O \setminus \{j,k\}}^T)^T \in \mathbb{R}^{2p-3}$ denote the parameters associated with the loss functions for the j th and k th observed

variables. The pairwise decorrelated score function for β_{jk} is defined as

$$S_{jk}(\beta_{j \vee k}) = \nabla_k \ell_j(\beta_{j, O \setminus j}) + \nabla_j \ell_k(\beta_{k, O \setminus k}) - (w_{jk}^*)^T \nabla_{-k} \ell_j(\beta_{j, O \setminus j}) - (w_{kj}^*)^T \nabla_{-j} \ell_k(\beta_{k, O \setminus k}). \quad (9)$$

The last two terms in (9) are constructed so that the effect of the nuisance parameters on the score function is asymptotically negligible (Section 3.2 of Ning et al., 2016). 200

The pairwise decorrelated score function (9) depends on the unknown quantities w_{jk}^* and w_{kj}^* . We estimate them using a Dantzig selector type estimator (Candès & Tao, 2007),

$$\hat{w}_{jk} = \arg \min_{w \in \mathbb{R}^{p-2}} \|w\|_1 \quad \text{subject to} \quad \left\| \nabla_{k, -k}^2 \ell_j(0, \hat{\beta}_{j, O \setminus \{j, k\}}) - w^T \nabla_{-k, -k}^2 \ell_j(0, \hat{\beta}_{j, O \setminus \{j, k\}}) \right\|_\infty \leq \lambda_w, \quad (10)$$

where $(0, \hat{\beta}_{j, O \setminus \{j, k\}})$ is an estimate of $\beta_{j, O \setminus j}$ obtained by solving (7) and replacing $\hat{\beta}_{jk}$ with zero, and λ_w is a non-negative tuning parameter. With some abuse of notation in (10), we use the notation $(0, \hat{\beta}_{j, O \setminus \{j, k\}})$ to indicate $(\hat{\beta}_{j1}, \dots, \hat{\beta}_{j, k-1}, 0, \hat{\beta}_{j, k+1}, \dots, \hat{\beta}_{jp})$. 205

The estimated pairwise decorrelated score function for testing $\beta_{jk}^* = 0$ is obtained by replacing $\beta_{j, O \setminus j}$, $\beta_{k, O \setminus k}$, w_{jk}^* , and w_{kj}^* in (9) with the estimated parameters $(0, \hat{\beta}_{j, O \setminus \{j, k\}})$, $(0, \hat{\beta}_{k, O \setminus \{j, k\}})$, \hat{w}_{jk} , and \hat{w}_{kj} , respectively, leading to

$$\hat{S}_{jk} = \nabla_k \ell_j(0, \hat{\beta}_{j, O \setminus \{j, k\}}) + \nabla_j \ell_k(0, \hat{\beta}_{k, O \setminus \{j, k\}}) - \hat{w}_{jk}^T \nabla_{-k} \ell_j(0, \hat{\beta}_{j, O \setminus \{j, k\}}) - \hat{w}_{kj}^T \nabla_{-j} \ell_k(0, \hat{\beta}_{k, O \setminus \{j, k\}}). \quad (11)$$

Let 210

$$\hat{\sigma}_{jk}^2 = \hat{\Sigma}_{jk, jk}^{jk} - 2\hat{\Sigma}_{jk, j \setminus k}^{jk} \hat{w}_{jk} - 2\hat{\Sigma}_{jk, k \setminus j}^{jk} \hat{w}_{kj} + \hat{w}_{jk}^T \hat{\Sigma}_{j \setminus k, j \setminus k}^{jk} \hat{w}_{jk} + \hat{w}_{kj}^T \hat{\Sigma}_{k \setminus j, k \setminus j}^{jk} \hat{w}_{kj}, \quad (12)$$

where $\hat{\Sigma}^{jk}$ is to be defined in (16). For a given significance level $0 < \alpha < 1$, our proposed pairwise decorrelated score test takes the form

$$\psi_{jk}(\alpha) = \begin{cases} 1, & |n^{1/2} \hat{S}_{jk} / \hat{\sigma}_{jk}| > \Phi^{-1}(1 - \alpha/2), \\ 0, & \text{otherwise,} \end{cases} \quad (13)$$

where $\Phi(x)$ is the standard normal cumulative distribution function. We will show in Section 3.3 that under the null hypothesis given in (8), the type I error of $\psi_{jk}(\alpha)$ converges to α . We summarise the overall procedure for conducting the pairwise decorrelated score test for (8) in Algorithm 1. 215

Algorithm 1. Pairwise decorrelated score test for testing $H_0 : \beta_{jk}^* = \beta_{kj}^* = 0$.

1. Obtain $\hat{\beta}_{j, O \setminus j}$ and $\hat{\beta}_{k, O \setminus k}$ by solving the optimization problem (7).
2. Obtain \hat{w}_{jk} and \hat{w}_{kj} from (10).
3. Calculate the estimated pairwise decorrelated score function \hat{S}_{jk} as in (11).
4. Calculate $\hat{\sigma}_{jk}^2$ as defined in (12).
5. Reject the null hypothesis $H_0 : \beta_{jk}^* = \beta_{kj}^* = 0$ if $|n^{1/2} \hat{S}_{jk} / \hat{\sigma}_{jk}| > \Phi^{-1}(1 - \alpha/2)$, where $0 < \alpha < 1$ is the given significance level.

3. THEORETICAL RESULTS

3.1. Notation

We use the Landau symbol $f(n) = \mathcal{O}\{g(n)\}$ to indicate the existence of a constant $C > 0$ such that $f(n) \leq Cg(n)$ for two sequences $f(n)$ and $g(n)$. We write $f(n) = \Omega\{g(n)\}$ to indicate $g(n) = \mathcal{O}\{f(n)\}$. In addition, we write $f(n) = o\{g(n)\}$ if $\lim_{n \rightarrow \infty} f(n)/g(n) \rightarrow 0$. We use the stochastic Landau symbol $f(n) = \mathcal{O}_{\mathbb{P}}\{g(n)\}$ to indicate that $f(n) = \mathcal{O}\{g(n)\}$ with high probability. For a vector $v = (v_1, \dots, v_d)^T \in \mathbb{R}^d$, we let $v^{\otimes 2}$ denote the outer product vv^T . For a symmetric matrix $M \in \mathbb{R}^{d \times d}$, we let $\|M\|_{\infty} = \max_{1 \leq j, j' \leq d} |M_{jj'}|$. Also, let $\Lambda_{\min}(M)$ and $\Lambda_{\max}(M)$ denote the minimum and maximum eigenvalues of M , respectively.

3.2. Parameter estimation

We provide an upper bound on the estimation error of $\hat{\beta}_{j, O \setminus j}$ obtained from solving (7). We study the asymptotic regime in which both n and p are allowed to grow, with R and h fixed. Proofs are deferred to the Supplementary Material. We first state an assumption on the first moment of the random variables and the local smoothness of the log-partition function.

Assumption 3. Let β_j^*, f_j^* be the true parameters in (3), and define the univariate function $\bar{A}_j(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ as

$$\bar{A}_j(u) = \log \left[\int \exp \left\{ ux_j + \frac{1}{2} \sum_{j=1}^{p+h} \sum_{k \neq j} \beta_{jk}^* x_j x_k + \sum_{j=1}^{p+h} f_j^*(x_j) \right\} d\nu(x) \right].$$

For all $j \in O$, we assume the following: (i) $|E(X_j)| \leq \kappa_m$, and (ii) $\max_{u: |u| \leq 1} \bar{A}_j''(u) \leq \kappa_h$.

Assumption 3 allows us to control the tail behaviors of the random variables. The same assumption has been used in recent work on the mixed graphical model (Chen et al., 2015).

Let $\mathcal{S}_j = \{k : \beta_{jk}^* \neq 0, k \in O \setminus j\}$ be the support set of $\beta_{j, O \setminus j}^*$ and let $s_j = |\mathcal{S}_j|$ be the cardinality of the set \mathcal{S}_j . Let $s_{\max} = \max_{j \in O} s_j$. Let κ_{\min}^2 , RE_{\min} , and $\rho_{q, \min}$ be the compatibility factor, restricted eigenvalue, and weak cone invertibility factor. These depend on the minimal eigenvalues of the Hessian matrix of the loss function, and will be defined rigorously in the Supplementary Material. These quantities are commonly used to establish upper bounds for estimation error in the context of ℓ_1 -penalised regression (Bickel et al., 2009; van de Geer & Bühlmann, 2009). We now establish upper bounds on the estimation error of $\hat{\beta}_{j, O \setminus j}$.

THEOREM 1. Let $\lambda = C(\log^5 p/n)^{1/2}$ for some constant $C > 0$. For $j \in O$, assume the event

$$\mathcal{A} = \left\{ \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left\| \begin{pmatrix} x_{ij}^r - x_{ij}^{r'} \end{pmatrix} \begin{pmatrix} x_{i, O \setminus j}^r - x_{i, O \setminus j}^{r'} \end{pmatrix} \right\|_{\infty} \leq M \right\}$$

and that $M s_{\max} \lambda / \kappa_{\min}^2 = o(1)$. Under Assumption 3, there exists a constant $C' > 0$ such that

$$\begin{aligned} \|\hat{\beta}_{j, O \setminus j} - \beta_{j, O \setminus j}^*\|_1 &\leq C' s_{\max} \lambda / \kappa_{\min}^2, \\ \|\hat{\beta}_{j, O \setminus j} - \beta_{j, O \setminus j}^*\|_2 &\leq C' (s_{\max})^{1/2} \lambda / \text{RE}_{\min}, \\ \|\hat{\beta}_{j, O \setminus j} - \beta_{j, O \setminus j}^*\|_q &\leq C' (s_{\max})^{1/q} \lambda / \rho_{q, \min}, \quad (q \geq 1), \end{aligned}$$

with probability at least $1 - p^{-1}$.

Theorem 1 generalises Theorem 4.4 in Yang et al. (arXiv:1412.8697). Interestingly, we obtain the same rate of convergence even when latent variables are present. Our rate of convergence

does not depend on the number of latent variables h . Theorem 1 holds with high probability conditioned on the event \mathcal{A} . For binary or categorical variables, \mathcal{A} holds with M constant. In the case of sub-exponential random variables, it can be shown that \mathcal{A} holds with high probability, with $M = C \log^2 p$ for a sufficiently large constant C . 250

The upper bounds on the estimation error in Theorem 1 depend on the quantities κ_{\min}^2 , RE_{\min} , and $\rho_{q,\min}$. These conditions can be bounded below by a positive constant when the random variables follow a multivariate Gaussian distribution.

THEOREM 2. Assume that $s_{\max}(\log^9 p/n)^{1/2} = o(1)$. Let 255

$$\{(X_{iO}^r)^T, X_{iH}^T\}^T \sim N(0, \Sigma), \quad \Sigma = \begin{pmatrix} \Sigma_{O,O} & \Sigma_{O,H} \\ \Sigma_{H,O} & \Sigma_{H,H} \end{pmatrix}.$$

Under Assumption 3, for n sufficiently large, the quantities κ_{\min}^2 , RE_{\min} , and $\rho_{q,\min}$ are larger than $C\Lambda_{\min}(\Sigma)$ with probability at least $1 - p^{-1}$ for some constant $C > 0$.

3.3. Pairwise decorrelated score test

In this section, we show that the type I error of the pairwise decorrelated score test in (13) converges to the desired significance level, under the null hypothesis $H_0 : \beta_{jk}^* = \beta_{kj}^* = 0$. We start by introducing some additional notation. Let 260

$$U_i^j(\beta_{j,O \setminus j}^*) = -\frac{2}{R(R-1)} \sum_{1 \leq r < r' \leq R} \frac{R_{ij}^{rr'}(\beta_{j,O \setminus j}^*)(x_{ij}^r - x_{ij}^{r'})(x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'})}{1 + R_{ij}^{rr'}(\beta_{j,O \setminus j}^*)} \in \mathbb{R}^{p-1}.$$

Furthermore, let $U_{ik}^j(\beta_{j,O \setminus j}^*)$ be the element in $U_i^j(\beta_{j,O \setminus j}^*)$ corresponding to the k th feature and let $U_{i,-k}^j(\beta_{j,O \setminus j}^*) \in \mathbb{R}^{p-2}$ be the vector obtained by removing the entry $U_{ik}^j(\beta_{j,O \setminus j}^*)$ from $U_i^j(\beta_{j,O \setminus j}^*)$. For $j, k \in O$, we let

$$g_i^{jk}(\beta_{j \vee k}^*) = \begin{Bmatrix} U_{ik}^j(\beta_{j,O \setminus j}^*) + U_{ij}^k(\beta_{k,O \setminus k}^*) \\ U_{i,-k}^j(\beta_{j,O \setminus j}^*) \\ U_{i,-j}^k(\beta_{k,O \setminus k}^*) \end{Bmatrix} \in \mathbb{R}^{2p-3} \quad (14)$$

and 265

$$\Sigma^{jk} = E \left[\left\{ g_i^{jk}(\beta_{j \vee k}^*) \right\}^{\otimes 2} \right] = \begin{Bmatrix} \Sigma_{jk,jk}^{jk} & \Sigma_{jk,j \setminus k}^{jk} & \Sigma_{jk,k \setminus j}^{jk} \\ (\Sigma_{jk,j \setminus k}^{jk})^T & \Sigma_{j \setminus k,j \setminus k}^{jk} & \Sigma_{j \setminus k,k \setminus j}^{jk} \\ (\Sigma_{jk,k \setminus j}^{jk})^T & (\Sigma_{j \setminus k,k \setminus j}^{jk})^T & \Sigma_{k \setminus j,k \setminus j}^{jk} \end{Bmatrix} \in \mathbb{R}^{(2p-3) \times (2p-3)}. \quad (15)$$

The quantity Σ^{jk} can be estimated using

$$\hat{\Sigma}^{jk} \left(0, \hat{\beta}_{j,O \setminus \{j,k\}}, \hat{\beta}_{k,O \setminus \{j,k\}} \right) = \frac{1}{n} \sum_{i=1}^n \left\{ g_i^{jk} \left(0, \hat{\beta}_{j,O \setminus \{j,k\}}, \hat{\beta}_{k,O \setminus \{j,k\}} \right) \right\}^{\otimes 2}. \quad (16)$$

In what follows, we write $\hat{\Sigma}^{jk}$ to indicate $\hat{\Sigma}^{jk}(0, \hat{\beta}_{j,O \setminus \{j,k\}}, \hat{\beta}_{k,O \setminus \{j,k\}})$.

We now state several assumptions. Recall from (9) that the pairwise decorrelated score function depends on the quantity $(w_{jk}^*)^T = (H_{k,-k}^j)^T (H_{-k,-k}^j)^{-1} \in \mathbb{R}^{p-2}$. The following assumption on the expected Hessian of the loss function guarantees that the pairwise decorrelated score function (9) is well-defined. 270

Assumption 4. Let $H^j = E\{\nabla^2 \ell_j(\beta_{j,O \setminus j}^*)\}$. For all $j \in O$, assume that

$$0 < \Lambda_{\text{lower}}^H \leq \Lambda_{\min}(H^j) \leq \Lambda_{\max}(H^j) \leq \Lambda_{\text{upper}}^H < \infty.$$

The next assumption guarantees that $g_i^{jk}(\beta_{j \vee k}^*)$ defined in (14) is not degenerate, in the sense that the variance of any linear combination of the elements of $g_i^{jk}(\beta_{j \vee k}^*)$ is not equal to zero. It is needed to guarantee the existence of the asymptotic variance of the score function (9).

Assumption 5. For $j, k \in O$, assume that $\Lambda_{\min}(\Sigma^{jk}) \geq \Lambda_{\text{lower}}^\Sigma > 0$.

The following theorem establishes that under the null hypothesis $H_0 : \beta_{jk}^* = \beta_{kj}^* = 0$, the type I error of $\psi_{jk}(\alpha)$ in (13) converges to α , and the associated p -value is asymptotically uniformly distributed in the $[0, 1]$ interval.

THEOREM 3. *Let the pairwise decorrelated score test with significance level $0 < \alpha < 1$, $\psi_{jk}(\alpha)$, be as defined in (13). We reject the null hypothesis $H_0 : \beta_{jk}^* = \beta_{kj}^* = 0$ if $\psi_{jk}(\alpha) = 1$. The associated p -value is defined as $\hat{p}_{jk} = 2\{1 - \Phi(|n^{1/2}\hat{S}_{jk}/\hat{\sigma}_{jk}|\})$, where*

$$\hat{\sigma}_{jk}^2 = \hat{\Sigma}_{jk,jk}^{jk} - 2\hat{\Sigma}_{jk,j \setminus k}^{jk} \hat{w}_{jk} - 2\hat{\Sigma}_{jk,k \setminus j}^{jk} \hat{w}_{kj} + \hat{w}_{jk}^\top \hat{\Sigma}_{j \setminus k, j \setminus k}^{jk} \hat{w}_{jk} + \hat{w}_{kj}^\top \hat{\Sigma}_{k \setminus j, k \setminus j}^{jk} \hat{w}_{kj}.$$

Under Assumptions 3–5 and scaling assumptions in Assumptions S1–S2 in the Supplementary Material, $\lim_{n \rightarrow \infty} \text{pr}\{\psi_{jk}(\alpha) = 1 \mid H_0\} = \alpha$ and \hat{p}_{jk} converges to a uniform distribution on the interval $[0, 1]$.

Results similar to Theorem 3 have been proven in the context of semiparametric regression and graphical models (Theorem 4.1 in Ning et al., 2016; Theorem 4.7 of Yang et al. (arXiv:1412.8697)).

4. SIMULATION STUDIES

4.1. Overview and competing proposals

Recall from Definition 1 that $\beta_{jk}^* \neq 0$ if and only if the j th and k th nodes are conditionally dependent, given all the other variables. To evaluate the performance across different methods, we define the true positive rate as the proportion of correctly identified non-zeros, and the false positive rate as the proportion of zeros that are incorrectly identified to be non-zeros. To examine the finite-sample performance of the pairwise decorrelated score test, we test the null hypothesis $H_0 : \beta_{jk}^* = 0$. The type I error and power are calculated as the proportion of falsely rejected H_0 and correctly rejected H_0 , respectively.

Five approaches are compared in our simulation studies: our proposal; the low-rank plus sparse latent variable Gaussian graphical model (Chandrasekaran et al., 2012); the semiparametric exponential family graphical model in Yang et al. (arXiv:1412.8697); the graphical lasso (Friedman et al., 2008); and the neighbourhood selection procedure (Meinshausen & Bühlmann, 2006; Ravikumar et al., 2011). Our proposal, Meinshausen & Bühlmann (2006), Ravikumar et al. (2011), and the semiparametric exponential family graphical model yield asymmetric estimates of the edge set. To symmetrise the edge set, we consider both the intersection and union rules described in Meinshausen & Bühlmann (2006), and report the best results for the competing proposals. We report our results using only the union rule.

Since the competing methods cannot accommodate replicates, we apply them to all nR observations, treating the replicates as independent samples. Our proposal, Friedman et al. (2008), Meinshausen & Bühlmann (2006), Ravikumar et al. (2011), and the semiparametric exponential

family graphical model each involves one tuning parameter. We applied a fine grid of tuning parameter values to obtain the curves shown in Figs. 1–4. There are two tuning parameters for Chandrasekaran et al. (2012). We set the second tuning parameter to equal ten times the first, and consider a fine grid of the first. Similar results were obtained for different ratios of the two tuning parameters.

4.2. Gaussian graphical models with latent variables

Let $\Theta = \Sigma^{-1}$ be the inverse covariance matrix of a Gaussian distribution, so that from Example 1, $\beta_{jk}^* = -\Theta_{jk}$. We partition Θ and Σ into

$$\Theta = \begin{pmatrix} \Theta_{O,O} & \Theta_{O,H} \\ \Theta_{H,O} & \Theta_{H,H} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{O,O} & \Sigma_{O,H} \\ \Sigma_{H,O} & \Sigma_{H,H} \end{pmatrix},$$

where $\Theta_{O,O}$, $\Theta_{O,H}$, and $\Theta_{H,H}$ encode the conditional dependence relationships among the observed variables, between the observed and latent variables, and among the latent variables. We construct $\Theta_{O,O}$ by randomly setting 10% of the off-diagonal entries to 0.3. For $\Theta_{O,H}$ and $\Theta_{H,H}$, we randomly set 80% of the off-diagonal entries to 0.3. To ensure the positive definiteness of Θ , we set $\Theta_{jj} = |\Lambda_{\min}(\Theta)| + 0.2$ for $j = 1, \dots, p + h$. Finally, we set $\Sigma = \Theta^{-1}$.

We first generate the latent variables x_{iH} for the n subjects from $N(0, \Sigma_{H,H})$. We then simulate the R replicates for each subject from the conditional distribution of the observed variables $N(\Sigma_{O,H} \Sigma_{H,H}^{-1} x_{iH}, \Sigma_{O,O} - \Sigma_{O,H} \Sigma_{H,H}^{-1} \Sigma_{H,O})$. The results for $n = 100$, $p = 100$, $h = \{2, 5, 10\}$, and $R = 10$, averaged over 100 datasets, are presented in Fig. 1.

In general, our proposal outperforms Friedman et al. (2008), Meinshausen & Bühlmann (2006), and the semiparametric exponential family graphical model, which do not model the latent variables. As shown in Fig. 1(a), our proposal has performance similar to Chandrasekaran et al. (2012), even though this is intended for the Gaussian setting which holds here, whereas our approach is semiparametric. As we increase the number of latent variables, the low-rank assumption of Chandrasekaran et al. (2012) is increasingly violated. Our proposal, which does not rely on the low-rank assumption, outperforms Chandrasekaran et al. (2012) when h is large.

Next, we investigate the role of the number of latent variables h and replicates R in the performance of our proposed method. We vary the ratio of R and h , while keeping $n = 100$ and $p = 100$ fixed. In addition, to study the tradeoff between n and R , we keep $p = 100$ and $h = 3$ fixed, and vary n and R with $nR = 600$. The results, averaged over 100 datasets, are shown in Fig. 2. From Figs. 2(a)–(b), we see that our proposal’s performance improves as we increase the ratio R/h . From Fig. 2(c), we see that the performance of our method improves when $R > 2$. This suggests that for a fixed experimental budget, that is, keeping nR fixed, it may be beneficial to collect more than two replicates per sample.

Our proposal relies on Assumption 1, which states that the latent variables are constant across replicates. We perform a sensitivity analysis by allowing the latent variables to vary across replicates within each subject. Let z_i^r be a h -vector with each element independently drawn from a uniform distribution $U[-\epsilon, \epsilon]$. We simulate the r th replicate for the i th observation from $N\{\Sigma_{O,H} \Sigma_{H,H}^{-1} (x_{iH} + z_i^r), \Sigma_{O,O} - \Sigma_{O,H} \Sigma_{H,H}^{-1} \Sigma_{H,O}\}$. We consider five values of $\epsilon = \{0, 1, 1.5, 2, 2.5\}$. Results averaged over 100 datasets are in Fig. 3, which shows that our proposal is robust to small perturbations of the latent variables.

We now perform the pairwise decorrelated score test described in Algorithm 1, in order to test the null hypothesis $H_0 : \beta_{jk}^* = 0$. The pairwise decorrelated score test involves two tuning parameters, λ in (7) and λ_w in (10). We select λ using 10-fold cross-validation, implemented in the R package glmnet. We use the R package fastclime to solve (10). We set $\lambda_w = 0.06$, so that

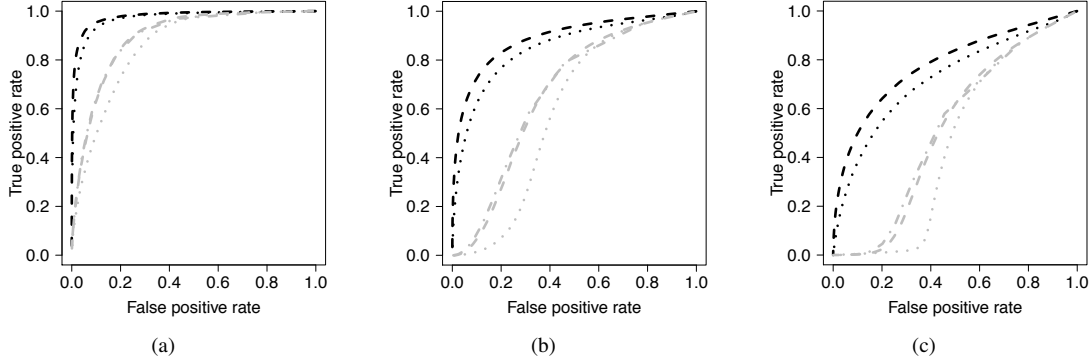


Fig. 1: Results for the simulation study for the Gaussian graphical model with $n = 100$, $p = 100$, and $R = 10$. Panels (a), (b), and (c) correspond to $h = \{2, 5, 10\}$ latent variables, respectively. The different curves represent our proposal (long-dashed), Chandrasekaran et al. (2012) (dots), Meinshausen & Bühlmann (2006) (grey long-dashed), Friedman et al. (2008) (grey dots), and the semiparametric exponential family graphical model in Yang et al. (arXiv:1412.8697) (grey dot-dashed).

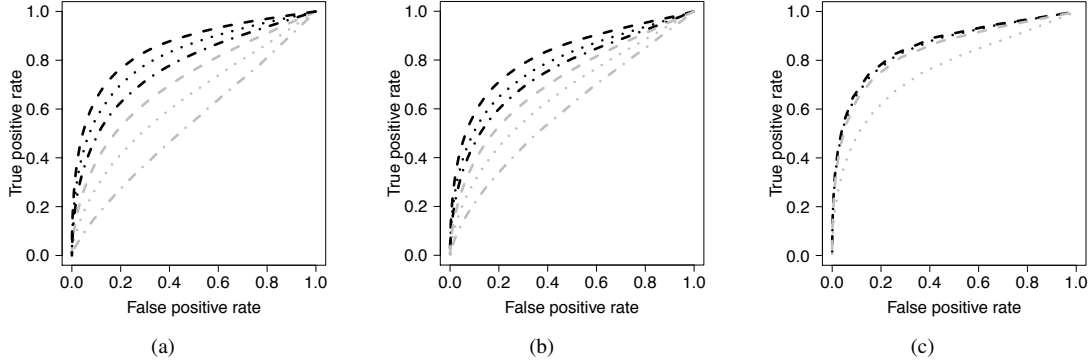


Fig. 2: Results of a simulation study investigating the relationship between h and R , and the tradeoff between n and R . Panels (a) and (b) are the results for $h = 8$ with $R = \{2, 4, 6, 8, 10, 12\}$, and $R = 6$ with $h = \{4, 5, 6, 8, 12, 24\}$, respectively, with $n = 100$ and $p = 100$. The curves represent different ratios R/h : 0.25 (grey dot-dashed), 0.5 (grey dots), 0.75 (grey long-dashed), 1 (dot-dashed), 1.25 (dots), and 1.5 (long-dashed). Panel (c) contains the results for $p = 100$ and $h = 3$, with different values of n and R such that $nR = 600$. The curves represent $n = 100$ and $R = 6$ (long-dashed), $n = 120$ and $R = 5$ (dots), $n = 150$ and $R = 4$ (grey long-dashed), and $n = 300$ and $R = 2$ (grey dots).

the estimates \hat{w}_{jk} and \hat{w}_{kj} contain a small number of non-zero entries. The results for $p = 100$, $R = 4$, and $h = 4$, over a range of sample sizes, are reported in Table 1. We see that the pairwise decorrelated score test is able to approximately control the type I error at level $\alpha = 0.05$.

4.3. Ising model with latent variables

We consider the Ising model with latent variables, as presented in Example 2. From Example 2, $\beta_{jk}^* = \Theta_{jk}$. We construct $\Theta_{O,O}$, $\Theta_{O,H}$, and $\Theta_{H,H}$ as in the previous section, but with nonzero

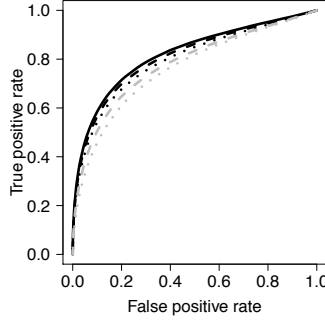


Fig. 3: Sensitivity analysis with $\text{Unif}[-\epsilon, \epsilon]$ noise added to each replicate, with $\epsilon = \{1, 1.5, 2, 2.5\}$. Results are for $n = 100$, $p = 100$, $h = 4$, and $R = 6$. The curves correspond to $\epsilon = 0$ (solid), $\epsilon = 1$ (long-dashed), $\epsilon = 1.5$ (dots), $\epsilon = 2$ (grey long-dashed), and $\epsilon = 2.5$ (grey dots).

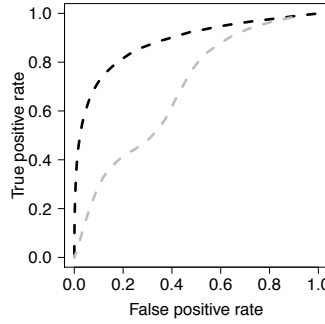


Fig. 4: Simulation results for the Ising model with $n = 100$, $p = 50$, $h = 5$, and $R = 10$, as described in Section 4.3. The curves represent our proposal (long-dashed) and the proposal of Ravikumar et al. (2011) (grey long-dashed).

entries drawn uniformly from $[-0.5, -0.25] \cup [0.25, 0.5]$. Furthermore, we do not require Θ to be positive definite. To obtain samples from the joint density (3), we employ a Gibbs sampler as described in Section 4 of Guo et al. (2015). The results for $n = 100$, $p = 50$, $h = 5$, and $R = 10$, averaged over 100 data sets, are presented in Fig. 4. Our proposal outperforms that of Ravikumar et al. (2011), which does not model the latent variables.

As in Section 4.2, we perform the pairwise decorrelated score test of the null hypothesis $H_0 : \beta_{jk}^* = 0$. We set $\lambda_w = 0.005$ in (10), so that the estimates \hat{w}_{jk} and \hat{w}_{kj} are sparse. The tuning parameter λ in (7) is again chosen by cross-validation. The type I error and power for $p = 50$, $R = 10$, and $h = 5$, over a range of sample sizes, are in Table 1. We see that the pairwise decorrelated score test is able to approximately control the type I error rate at level $\alpha = 0.05$.

5. APPLICATION TO ADHD-200 DATA

We applied our method to the ADHD-200 data (Biswal et al., 2010). The data consist of resting state functional magnetic resonance images on 197 subjects who have been diagnosed with attention deficit hyperactivity disorder, and 491 control subjects. The number of images for each subject ranges from 76 to 276. Covariates such as age, gender, site, and intelligence quotient

Table 1: Type I error and power of the pairwise decorrelated score test at the 5% significance level are calculated as the % of falsely rejected and correctly rejected null hypotheses, respectively, over 2000 data sets. Data were generated under the Gaussian graphical model with latent variables with $p = 100$, $R = 4$, and $h = 4$. Data were generated under the Ising model with latent variables with $p = 50$, $R = 10$, and $h = 5$

		$n = 50$	$n = 100$	$n = 200$	$n = 300$	$n = 400$
Gaussian	Type I error	9	7	6	5	5
	Power	18	27	45	59	70
Ising	Type I error	7	5	5	5	5
	Power	30	46	73	87	95

are also available. Similar to Power et al. (2011) and Qiu et al. (2016), we use 264 seed regions of interest to define the nodes in the graphical model.

We treat the images for each subject as replicates, and treat the covariates such as age and gender as latent variables. However, certain covariates such as age and gender serve as confounders that may alter the conditional dependence relationships among the variables. For instance, Qiu et al. (2016) showed that the brain networks at ages 7, 12, and 22 years are quite different. Biswal et al. (2010) showed that males and females have different brain connectivity networks. Thus, standard techniques for estimating graphical models that do not model the latent variables may yield inaccurate network estimates.

After removing subjects with missing values, we consider 465 control subjects in the dataset. For computational purposes, we choose $R = 10$ replicates randomly for each subject. Assumption 2 may not hold, since the replicate brain images for a given subject are very likely to be dependent. We standardize each seed region to have mean zero and standard deviation one for each subject. Our proposal (7) involves one tuning parameter λ . For visualization, we set $\lambda = 0.2$ so that the estimated network is sparse, but in practice, λ can be chosen by cross-validation. We then symmetrise our estimates using the intersection rule described in Section 4. This yields an estimated network with 376 edges. Figs. 5(a)–(c) show coronal, sagittal, and transverse snapshots of the estimated brain connectivity network.

We compare our proposal to that of Friedman et al. (2008), which does not model the latent variables. We perform their proposal by treating the replicates as independent observations. For ease of comparison, the tuning parameter for Friedman et al. (2008) is chosen to yield 376 edges. The coronal, sagittal, and transverse snapshots of the estimated brain connectivity network from Friedman et al. (2008) are plotted in Figs. 5(d)–(f).

The two estimated networks are somewhat different. For instance, we see from Figs. 5(b) and 5(e) that the lower region of the brain connectivity network estimated by their proposal is more densely connected than that of our proposal. This might be a consequence of marginalizing over the latent variables, as discussed in Chandrasekaran et al. (2012). In contrast, edges in the network estimated by our proposal seem to be more spread throughout the network.

6. DISCUSSION

Our proposal can be generalised beyond estimating latent variable graphical models. For instance, in the context of regression, unmeasured confounders may remain constant across replicates. Without adjusting for these confounders, it can be shown that the estimated regression coefficients for the observed variables are biased. Using the ideas laid out in this paper, one can estimate the parameter of interest accurately by treating the confounders as nuisance parameters.

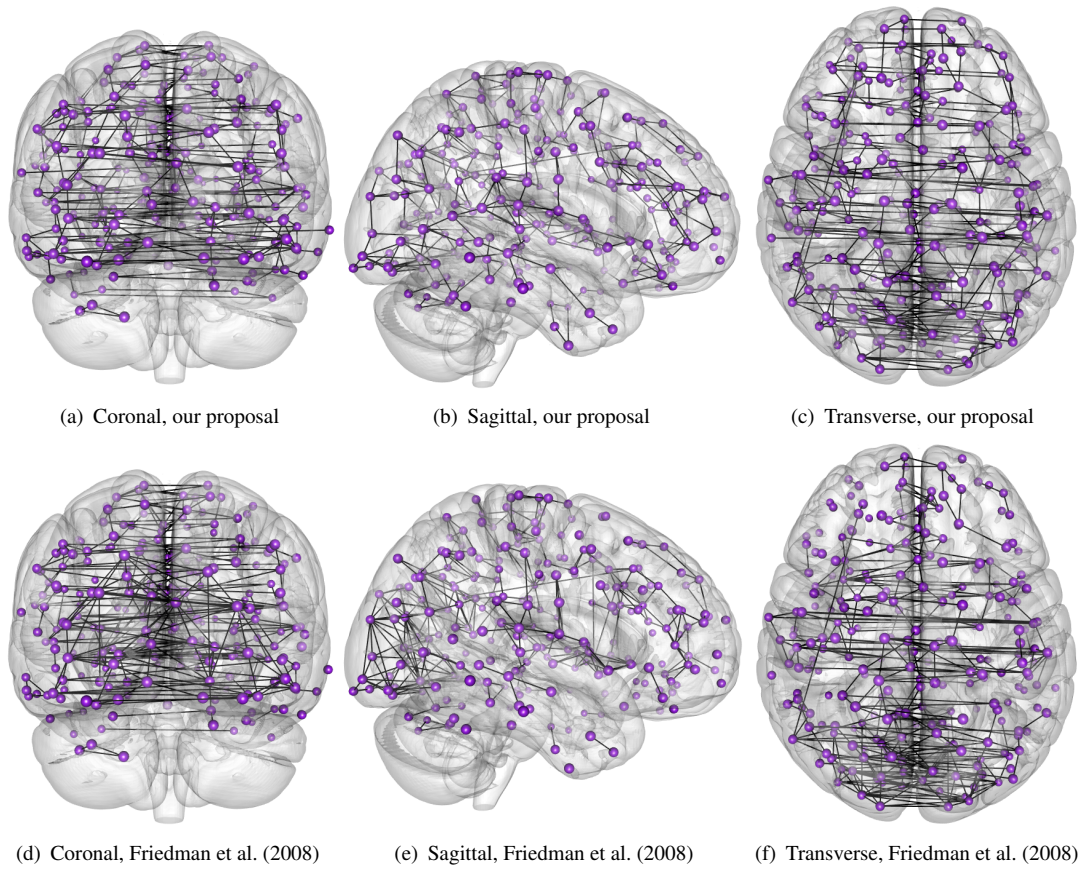


Fig. 5: Coronal, sagittal, and transverse snapshots of the estimated brain connectivity networks resulting from our proposal and Friedman et al. (2008). Panels (a)–(c) and panels (d)–(f) contain the estimated networks from our proposal and Friedman et al. (2008), respectively.

Our model requires that the replicates are mutually conditionally independent given the latent variables; this is laid out in Assumption 2. In future work, it would be interesting to study whether that assumption can be relaxed.

410

An R package `latentGraph` will be made available on CRAN.

ACKNOWLEDGEMENT

We thank the editor, an associate editor, and two reviewers for helpful comments, Shizhe Chen for responding to our inquiries, and Huitong Qiu for providing us R code to plot Fig. 5.

SUPPLEMENTARY MATERIAL

415

Supplementary material available at *Biometrika* online includes proofs of the theoretical results and the scaling assumptions used in Theorem 3.

REFERENCES

- BICKEL, P. J., RITOV, Y. & TSYBAKOV, A. B. (2009). Simultaneous analysis of Lasso and Dantzig selector. *Ann. Statist.* **37**, 1705–1732.
- BISWAL, B. B., MENNES, M., ZUO, X.-N. et al. (2010). Toward discovery science of human brain function. *Proc. Natl. Acad. Sci.* **107**, 4734–4739.
- BOX, G. E., HUNTER, J. S. & HUNTER, W. G. (2005). *Statistics for Experimenters: Design, Innovation, and Discovery*. New York: Wiley-Interscience, 2nd ed.
- CAI, T. T., LIU, W. & LUO, X. (2011). A constrained ℓ_1 minimization approach to sparse precision matrix estimation. *J. Am. Statist. Assoc.* **106**, 594–607.
- CANDÈS, E. J. & TAO, T. (2007). The Dantzig selector: statistical estimation when p is much larger than n . *Ann. Statist.* **35**, 2313–2351.
- CHANDRASEKARAN, V., PARRILO, P. A. & WILLSKY, A. S. (2012). Latent variable graphical model selection via convex optimization. *Ann. Statist.* **40**, 1935–1967.
- CHEN, S., WITTEN, D. M. & SHOJAIE, A. (2015). Selection and estimation for mixed graphical models. *Biometrika* **102**, 47–64.
- FELLINGHAUER, B., BÜHLMANN, P., RYFFEL, M., VON RHEIN, M. & REINHARDT, J. D. (2013). Stable graphical model estimation with random forests for discrete, continuous, and mixed variables. *Comp. Statist. Data Anal.* **64**, 132–152.
- FRIEDMAN, J. H., HASTIE, T. J. & TIBSHIRANI, R. J. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**, 432–441.
- GUO, J., CHENG, J., LEVINA, E., MICHAILIDIS, G. & ZHU, J. (2015). Estimating heterogeneous graphical models for discrete data with an application to roll call voting. *Ann. Appl. Statist.* **9**, 821–848.
- LEE, J. D. & HASTIE, T. J. (2015). Learning the structure of mixed graphical models. *J. Comp. Graph. Statist.* **25**, 230–253.
- LIU, H., HAN, F., YUAN, M., LAFFERTY, J. D. & WASSERMAN, L. (2012). High-dimensional semiparametric Gaussian copula graphical models. *Ann. Statist.* **40**, 2293–2326.
- LIU, H., LAFFERTY, J. D. & WASSERMAN, L. (2009). The nonparanormal: semiparametric estimation of high dimensional undirected graphs. *J. Mach. Learn. Res.* **10**, 2295–2328.
- MEINSHAUSEN, N. & BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.* **34**, 1436–1462.
- MONTGOMERY, D. C. (2008). *Design and Analysis of Experiments*. New York: John Wiley & Sons, 8th ed.
- NING, Y. & LIU, H. (2016). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Ann. Statist.*, in press.
- NING, Y., ZHAO, T. & LIU, H. (2016). A likelihood ratio framework for high dimensional semiparametric regression. *Ann. Statist.*, in press.
- PENG, J., WANG, P., ZHOU, N. & ZHU, J. (2009). Partial correlation estimation by joint sparse regression models. *J. Am. Statist. Assoc.* **104**, 735–746.
- POWER, J. D., COHEN, A. L., NELSON, S. M. et al. (2011). Functional network organization of the human brain. *Neuron* **72**, 665–678.
- QIU, H., HAN, F., LIU, H. & CAFFO, B. (2016). Joint estimation of multiple graphical models from high dimensional time series. *J. R. Statist. Soc. B* **78**, 487–504.
- RAVIKUMAR, P., WAINWRIGHT, M. J. & LAFFERTY, J. D. (2010). High-dimensional Ising model selection using ℓ_1 -regularized logistic regression. *Ann. Statist.* **38**, 1287–1319.
- RAVIKUMAR, P., WAINWRIGHT, M. J., RASKUTTI, G. & YU, B. (2011). High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electron. J. Statist.* **5**, 935–980.
- ROTHMAN, A. J., BICKEL, P. J., LEVINA, E. & ZHU, J. (2008). Sparse permutation invariant covariance estimation. *Electron. J. Statist.* **2**, 494–515.
- SUN, T. & ZHANG, C.-H. (2013). Sparse matrix inversion with scaled lasso. *J. Mach. Learn. Res.* **14**, 3385–3418.
- VAN DE GEER, S. A. & BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the lasso. *Electron. J. Statist.* **3**, 1360–1392.
- VOORMAN, A., SHOJAIE, A. & WITTEN, D. M. (2014). Graph estimation with joint additive models. *Biometrika* **101**, 85–101.
- YANG, E., RAVIKUMAR, P., ALLEN, G. I. & LIU, Z. (2015). Graphical models via univariate exponential family distributions. *J. Mach. Learn. Res.* **16**, 3813–3847.
- YUAN, M. & LIN, Y. (2007). Model selection and estimation in the Gaussian graphical model. *Biometrika* **94**, 19–35.

[Received MY . Revised MY]