*Biometrika* (2016), **99**, 1, *pp*. 1–17 © 2016 Biometrika Trust *Printed in Great Britain* 

# 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

#### 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

# AND HAN LIU

Department of Operations Research and Financial Engineering, Princeton University, Princeton, New Jersey 08544, U.S.A. hanliu@princeton.edu

#### 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.<sup>20</sup> 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.

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

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

30

5

10

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.,

<sup>45</sup> 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
 <sup>50</sup> 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

<sup>55</sup> 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.

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  $X_{-j} = (X_1, \ldots, X_{j-1}, X_{j+1}, \ldots, X_p)^T$ . 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

$$p(x_j \mid x_{-j}) = \exp\left\{x_j \beta_{j,-j}^{\mathrm{T}} x_{-j} + f_j(x_j) - A_j(\beta_j, f_j)\right\},\tag{1}$$

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  $x_{\cdot j} = (x_{1j}, \ldots, x_{nj})^{\mathrm{T}}$ , and let  $x_j^{(\cdot)}$  and  $z_{\cdot j}$  be the order and rank statistics of  $x_{\cdot j}$ , respec-

2

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

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

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  $X = (X_O^T, X_H^T)^T$  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  $O = \{1, \ldots, p\}$  and  $H = \{p+1, \ldots, p+h\}$  denote the index sets of the observed and latent random variables, respectively.

DEFINITION 1. A (p+h)-dimensional random vector  $X = (X_O^T, X_H^T)^T$  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

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

where  $f_i(x_i)$  is some possibly unknown function and  $A_i(\beta_i, f_i)$  is the log-partition function.

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

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

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

95

100

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

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

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

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

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  $R_1 = \cdots = R_n = R$ , 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  $i = 1, \ldots, n$  and  $r = 1, \ldots, R$ . 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  $1 \le r' \le r \le R$ .

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

<sup>140</sup> 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.

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

$$p\left(X_{ij}^{r} = x_{ij}^{r}, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^{r}, x_{i,-j}^{r'}\right) = p\left\{X_{ij}^{r} = x_{ij}^{r}, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^{r}, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\right\} p\left\{x_{ij}^{(r,r')} \mid x_{i,-j}^{r}, x_{i,-j}^{r'}\right\}.$$
(4)

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

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

$$p\left\{X_{ij}^{r} = x_{ij}^{r}, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,-j}^{r}, x_{i,-j}^{r'}, x_{ij}^{(r,r')}\right\} = \left\{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})\right\}^{-1},$$
(5)  
where  $R_{ij}^{rr'}(\beta_{j,O\setminus j}) = \exp\{-(x_{ij}^{r} - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})\}.$ 

In the absence of latent variables, similar results were established in the context of the semiparametric 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 <sup>160</sup> which applied the approach outlined in Section 2.1.

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

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

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  $\beta_{j,O\setminus j}$ .

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,

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

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

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

155

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

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

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

where  $\lambda$  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  $\tilde{x}_{ij}^{rr'} = \operatorname{sign}(x_{ij}^r - x_{ij}^{r'})$ 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  $\ell_1$ -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  $\beta_{i,O\setminus i}^*$  and  $\beta_{k,O\setminus k}^*$ , that is,

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

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

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

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

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

$$\begin{split} H^{j}_{k,-k} &= E\left\{\nabla^{2}_{k,-k}\ell_{j}(\beta^{*}_{j,O\setminus j})\right\} \in \mathbb{R}^{p-2}, \qquad H^{j}_{-k,-k} = E\left\{\nabla^{2}_{-k,-k}\ell_{j}(\beta^{*}_{j,O\setminus j})\right\} \in \mathbb{R}^{(p-2)\times(p-2)}.\\ \text{Let } (w^{*}_{jk})^{\mathrm{T}} &= (H^{j}_{k,-k})^{\mathrm{T}}(H^{j}_{-k,-k})^{-1} \in \mathbb{R}^{p-2}, \text{ and let } \beta_{j\vee k} = (\beta_{jk},\beta^{\mathrm{T}}_{j,O\setminus\{j,k\}},\beta^{\mathrm{T}}_{k,O\setminus\{j,k\}})^{\mathrm{T}} \in \mathbb{R}^{2p-3} \text{ denote the parameters associated with the loss functions for the jth and kth observed expression. \end{split}$$

variables. The pairwise decorrelated score function for  $\beta_{ik}$  is defined as

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

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

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),

$$\hat{w}_{jk} = \underset{w \in \mathbb{R}^{p-2}}{\operatorname{arg\,min}} \|w\|_1 \quad \text{subject to} \quad \left\|\nabla_{k,-k}^2 \ell_j(0,\hat{\beta}_{j,O\setminus\{j,k\}}) - w^{\mathrm{T}} \nabla_{-k,-k}^2 \ell_j(0,\hat{\beta}_{j,O\setminus\{j,k\}})\right\|_{\infty} \leq \lambda_w,$$
(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  $\lambda_w$  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  $(\hat{\beta}_{j1}, \ldots, \hat{\beta}_{j,k-1}, 0, \hat{\beta}_{j,k+1}, \ldots, \hat{\beta}_{jp})$ . The estimated pairwise decorrelated score function for testing  $\beta_{jk}^* = 0$  is obtained by replacing

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

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

Let

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

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

$$\psi_{jk}(\alpha) = \begin{cases} 1, & \left| n^{1/2} \hat{S}_{jk} / \hat{\sigma}_{jk} \right| > \Phi^{-1}(1 - \alpha/2), \\ 0, & \text{otherwise}, \end{cases}$$
(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  $\alpha$ . We summarise the overall procedure for conducting the pairwise decorrelated score test for (8) in Algorithm 1.

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  $\hat{\sigma}_{ik}^2$  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.

#### K. M. TAN, Y. NING, D. M. WITTEN AND H. LIU

#### 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 > 0such 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  $\lim_{n\to\infty} f(n)/g(n) \to 0$ . 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  $v = (v_1, \ldots, v_d)^{\mathrm{T}} \in \mathbb{R}^d$ , we let  $v^{\otimes 2}$  denote the outer product  $vv^{\mathrm{T}}$ . 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  $\beta_j^*, f_j^*$  be the true parameters in (3), and define the univariate function  $\bar{A}_i(\cdot) : \mathbb{R} \to \mathbb{R}$  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)  $\max_{u:|u|\leq 1} \overline{A}''_j(u) \leq \kappa_h$ .

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  $S_j = \{k : \beta_{jk}^* \neq 0, k \in O \setminus j\}$  be the support set of  $\beta_{j,O\setminus j}^*$  and let  $s_j = |S_j|$  be the cardinality of the set  $S_j$ . Let  $s_{\max} = \max_{j \in O} s_j$ . Let  $\kappa_{\min}^2$ , RE<sub>min</sub>, 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  $\ell_1$ -penalised regression (Bickel et al., 2009; van de Geer & Bühlmann,

2009). We now establish upper bounds on the estimation error of  $\hat{\beta}_{i,O\setminus j}$ .

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

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

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

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

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.

The upper bounds on the estimation error in Theorem 1 depend on the quantities  $\kappa_{\min}^2$ , RE<sub>min</sub>, 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

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

Under Assumption 3, for n sufficiently large, the quantities  $\kappa_{\min}^2$ , RE<sub>min</sub>, 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

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

Furthermore, let  $U_{ik}^{j}(\beta_{i,O\setminus j}^{*})$  be the element in  $U_{i}^{j}(\beta_{i,O\setminus j}^{*})$  corresponding to the kth 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  $U_i^j(\beta_{i|O\setminus i}^*)$ . For  $j,k \in O$ , we let

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

and

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

The quantity  $\Sigma^{jk}$  can be estimated using

$$\hat{\Sigma}^{jk}\left(0,\hat{\beta}_{j,O\setminus\{j,k\}},\hat{\beta}_{k,O\setminus\{j,k\}}\right) = \frac{1}{n}\sum_{i=1}^{n}\left\{g_{i}^{jk}(0,\hat{\beta}_{j,O\setminus\{j,k\}},\hat{\beta}_{k,O\setminus\{j,k\}})\right\}^{\otimes 2}.$$
(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  $(w_{jk}^*)^{\mathrm{T}} = (H_{k,-k}^j)^{\mathrm{T}}(H_{-k,-k}^j)^{-1} \in \mathbb{R}^{p-2}$ . The following assumption on the expected Hessian of the loss function guarantees that the pairwise decorrelated score 270 function (9) is well-defined.

255

265

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

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

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  $\Lambda_{\min}(\Sigma^{jk}) \ge \Lambda_{\text{lower}}^{\Sigma} > 0$ .

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  $\alpha$ , 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

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

Under Assumptions 3–5 and scaling assumptions in Assumptions S1–S2 in the Supplementary Material,  $\lim_{n\to\infty} \operatorname{pr}\{\psi_{jk}(\alpha) = 1 \mid H_0\} = \alpha$  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  $H_0$ 

295

290

285

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 310 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  $\Theta$  and  $\Sigma$  into

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

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  $\Theta_{H,H}$ , 320 we randomly set 80% of the off-diagonal entries to 0.3. To ensure the positive definiteness of  $\Theta$ , we set  $\Theta_{jj} = |\Lambda_{\min}(\Theta)| + 0.2$  for j = 1, ..., 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 = 100, p =325  $\{2, 5, 10\}$ , 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 335 p = 100 fixed. In addition, to study the tradeoff between n and R, we keep p = 100 and h = 3fixed, 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 340 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 rth replicate for the *i*th obser-345 vation 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  $H_0: \beta_{ik}^* = 0$ . The pairwise decorrelated score test involves two tuning 350 parameters,  $\lambda$  in (7) and  $\lambda_w$  in (10). We select  $\lambda$  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

330

315



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  $\Theta$  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 uning parameter  $\lambda$  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. 375

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

380

385

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  $\lambda$ . For visualization, we set  $\lambda = 0.2$ so that the estimated network is sparse, but in practice,  $\lambda$  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

390

395

400

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).

yield inaccurate network estimates.

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

405

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.

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

410

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<sub>1</sub> 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 heterogenous graphical models for discrete data with an application to roll call voting. *Ann. Appl. Statist.* 9, 821–848.
- 440 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  $\ell_{1}$ -regularized logistic regression. *Ann. Statist.* **38**, 1287–1319.
  - RAVIKUMAR, P., WAINWRIGHT, M. J., RASKUTTI, G. & YU, B. (2011). High-dimensional covariance estimation by minimizing  $\ell_1$ -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.
- 465 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.
- 470 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]

Biometrika (2016), xx, x, pp. 1-28 Printed in Great Britain

# Supplementary material for '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

10

15

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

#### AND HAN LIU

Department of Operations Research and Financial Engineering, Princeton University, Princeton, New Jersey 08544, U.S.A. hanliu@princeton.edu

#### S1. LIST OF NOTATION

In this section, we define some notation that will be used throughout this supplementary material. Recall from § 3.2 in the main paper that  $\beta_{j,O\setminus j}^*$  is the vector of true parameter values of interest,  $S_j$  is the support set of  $\beta_{j,O\setminus j}^*$ , and  $s_j = |S_j|$ . Let  $s_{\max} = \max_{j \in O} s_j$ . Let 20

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

The gradient of  $\ell_j(\beta_{i,O\setminus j})$  is

$$\nabla \ell_{j}(\beta_{j,O\setminus j}) = -\frac{1}{n} \sum_{i=1}^{n} \left\{ \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} \frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(x_{ij}^{r} - x_{ij}^{r'})(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})}{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})} \right\}$$
$$= \frac{1}{n} \sum_{i=1}^{n} U_{i}^{j}(\beta_{j,O\setminus j}) \in \mathbb{R}^{p-1},$$
(S1)

where  $U_1^j(\beta_{j,O\setminus j}), \ldots, U_n^j(\beta_{j,O\setminus j})$  are independent and identically distributed random variables. The Hessian of  $\ell_j(\beta_{j,O\setminus j})$  is

$$\nabla^{2} \ell_{j}(\beta_{j,O\setminus j}) = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} \frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(x_{ij}^{r} - x_{ij}^{r'})^{2}(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})^{\otimes 2}}{\{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})\}^{2}} \right]$$
  
$$\in \mathbb{R}^{(p-1) \times (p-1)}.$$

For notational convenience, we write

$$h_{ij}^{rr'}(\beta_{j,O\setminus j}) = -\frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(x_{ij}^r - x_{ij}^{r'})(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})}{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})} \in \mathbb{R}^{p-1}$$
(S2)

and, for  $k \in O \setminus j$ ,

$$h_{ijk}^{rr'}(\beta_{j,O\setminus j}) = -\frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(x_{ij}^r - x_{ij}^{r'})(x_{ik}^r - x_{ik}^{r'})}{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})} \in \mathbb{R}.$$
(S3)

Similarly, we let

$$T_{ij}^{rr'}(\beta_{j,O\backslash j}) = \frac{R_{ij}^{rr'}(\beta_{j,O\backslash j})(x_{ij}^r - x_{ij}^{r'})^2(x_{i,O\backslash j}^r - x_{i,O\backslash j}^{r'})^{\otimes 2}}{\{1 + R_{ij}^{rr'}(\beta_{j,O\backslash j})\}^2} \in \mathbb{R}^{(p-1)\times(p-1)}$$
(S4)

and, for  $k, l \in O \setminus j$ ,

$$T_{ijkl}^{rr'}(\beta_{j,O\backslash j}) = \frac{R_{ij}^{rr'}(\beta_{j,O\backslash j})(x_{ij}^r - x_{ij}^{r'})^2(x_{ik}^r - x_{ik}^{r'})(x_{il}^r - x_{il}^{r'})}{\{1 + R_{ij}^{rr'}(\beta_{j,O\backslash j})\}^2} \in \mathbb{R}.$$
(S5)

Therefore, the gradient and Hessian of the loss function can also be written as

$$\nabla \ell_j(\beta_{j,O\setminus j}) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} h_{ij}^{rr'}(\beta_{j,O\setminus j}) \right\} \in \mathbb{R}^{p-1}$$
(S6)

and

$$\nabla^2 \ell_j(\beta_{j,O\setminus j}) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} T_{ij}^{rr'}(\beta_{j,O\setminus j}) \right\} \in \mathbb{R}^{(p-1) \times (p-1)}.$$
 (S7)

35

#### S2. SCALING ASSUMPTIONS IN THEOREM 3

We state two assumptions on the scaling of n and p and on the magnitude of the regularization parameters  $\lambda$  and  $\lambda_w$  in (7) and (10), respectively. The following assumption is needed to show the asymptotic normality of (9).

Assumption S1. Let M be as defined in Theorem 1 and let  $w_0 = \max_{j,k\in O} \|w_{jk}^*\|_1$ . Furthermore, let  $s'_{jk} = \|w_{jk}^*\|_0$  and  $s'_{\max} = \max_{j,k\in O} s'_{jk}$ . Assume that  $Ms_{\max}\lambda/\kappa_{\min}^2 = o(1)$ ,  $s'_{\max}\lambda_w = o(1)$ ,

$$\lambda_w = \Omega \left\{ w_0 \left( \frac{M s_{\max} \lambda}{\kappa_{\min}^2} + \lambda \log^2 p \right) \right\}, \qquad \lim_{n \to \infty} n^{1/2} \left( \frac{s_{\max} \lambda_w \lambda}{\kappa_{\min}^2} \right) = 0$$

and

$$\lim_{n \to \infty} n^{1/2} s'_{\max} \lambda_w \lambda = 0.$$

Next, we state an additional scaling assumption to guarantee that  $\hat{\sigma}_{jk}$  in (12) is a consistent estimator of its asymptotic variance.

Assumption S2. For any  $j, k \in O$ , assume that

$$(1+w_0+w_0^2)\left(\frac{s_{\max}\lambda\log^6 p}{\kappa_{\min}^2}\right) = o(1), \qquad w_0 s'_{\max}\lambda_w = o(1)$$

# S3. Proof of the results in $\S\,2{\cdot}2$

S3.1. *Proof of Proposition* 1

Under Assumption 1, we obtain

$$\begin{aligned} & \operatorname{pr} \left\{ X_{ij}^{r} = x_{ij}^{r}, X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r}, x_{i,O\setminus j}^{r'}, x_{iH}, x_{ij}^{(r,r')} \right\} \\ &= \operatorname{pr} (X_{ij}^{r} = x_{ij}^{r} \mid x_{i,O\setminus j}^{r}, x_{iH}) \operatorname{pr} (X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r'}, x_{iH}) \Big/ \\ & \left\{ \operatorname{pr} (X_{ij}^{r} = x_{ij}^{r} \mid x_{i,O\setminus j}^{r}, x_{iH}) \operatorname{pr} (X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r'}, x_{iH}) \right. \\ & \left. + \operatorname{pr} (X_{ij}^{r} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r}, x_{iH}) \operatorname{pr} (X_{ij}^{r'} = x_{ij}^{r} \mid x_{i,O\setminus j}^{r'}, x_{iH}) \right\} \\ &= \left\{ 1 + \frac{\operatorname{pr} (X_{ij}^{r} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r}, x_{iH}) \operatorname{pr} (X_{ij}^{r'} = x_{ij}^{r} \mid x_{i,O\setminus j}^{r'}, x_{iH})}{\operatorname{pr} (X_{ij}^{r} = x_{ij}^{r} \mid x_{i,O\setminus j}^{r}, x_{iH}) \operatorname{pr} (X_{ij}^{r'} = x_{ij}^{r'} \mid x_{i,O\setminus j}^{r'}, x_{iH})} \right\}^{-1} \\ &= \left[ 1 + \exp\left\{ - (x_{ij}^{r} - x_{ij}^{r'}) \beta_{j,O\setminus j}^{\mathrm{T}} (x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'}) \right\} \right]^{-1}, \end{aligned}$$

where the first equality follows from Assumption 2, that the replicates are mutually independent <sup>55</sup> given the latent variables.

# S3-2. Proof of Proposition 2

Recall from § 2·3 that  $x_{ij}^{(r,r')} = \{\min(x_{ij}^r, x_{ij}^{r'}), \max(x_{ij}^r, x_{ij}^{r'})\}$  are the order statistics of a pair of replicates for the *i*th subject. Let  $X_{ij}^{(r,r')} = \{\min(X_{ij}^r, X_{ij}^{r'}), \max(X_{ij}^r, X_{ij}^{r'})\}$ . For notational convenience, we let  $\mathcal{B}_{ij}^{rr'}$  denote the event  $\{X_{i,-j}^r = x_{i,-j}^r, X_{i,-j}^{r'} = x_{i,-j}^{r'}, X_{ij}^{(r,r')} = x_{ij}^{(r,r')}\}$ . By Proposition 1, one can see that the conditional distribution of  $X_{ij}^r$  and  $X_{ij}^{r'}$  given  $\mathcal{B}_{ij}^{rr'}$  is binomial with

$$\operatorname{pr}\left(X_{ij}^{r} = x_{ij}^{r}, X_{ij}^{r'} = x_{ij}^{r'} \mid \mathcal{B}_{ij}^{rr'}\right) = \frac{1}{1 + \exp\left\{-(x_{ij}^{r} - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})\right\}}$$
(S8)

and

$$\operatorname{pr}\left(X_{ij}^{r} = x_{ij}^{r'}, X_{ij}^{r'} = x_{ij}^{r} \mid \mathcal{B}_{ij}^{rr'}\right) = \frac{\exp\left\{-(x_{ij}^{r} - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})\right\}}{1 + \exp\left\{-(x_{ij}^{r} - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})\right\}}.$$
 (S9)

45

Recall from (S2) that

$$\begin{split} h_{ij}^{rr'}(\beta_{j,O\setminus j}) &= -\frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(x_{ij}^r - x_{ij}^{r'})(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})}{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})} \\ &= -\frac{\exp\{-(x_{ij}^r - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})\}(x_{ij}^r - x_{ij}^{r'})(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})}{1 + \exp\{-(x_{ij}^r - x_{ij}^{r'})\beta_{j,O\setminus j}^{\mathrm{T}}(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})\}}. \end{split}$$

<sup>65</sup> The conditional expectation of  $h_{ij}^{rr'}(\beta^*_{j,O\setminus j})$  given  $\mathcal{B}_{ij}^{rr'}$  takes the form

$$\begin{split} & E\{h_{ij}^{rr'}(\beta_{j,O\setminus j}^{*}) \mid \mathcal{B}_{ij}^{rr'}\} \\ &= E\left[-\frac{\exp\{-(X_{ij}^{r}-X_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}(X_{ij}^{r}-X_{ij}^{r'})(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})}{1+\exp\{-(X_{ij}^{r}-X_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}}\right| \mathcal{B}_{ij}^{rr'}\right] \\ &= -\frac{\exp\{-(x_{ij}^{r}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}(x_{ij}^{r}-x_{ij}^{r'})(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})}{1+\exp\{-(x_{ij}^{r}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}} \operatorname{pr}(X_{ij}^{r}=x_{ij}^{r},X_{ij}^{r'}=x_{ij}^{r'}\mid \mathcal{B}_{ij}^{rr'}) \\ &-\frac{\exp\{-(x_{ij}^{r'}-x_{ij}^{r})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})}{1+\exp\{-(x_{ij}^{r'}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}} \operatorname{pr}(X_{ij}^{r}=x_{ij}^{r'},X_{ij}^{r'}=x_{ij}^{r'}\mid \mathcal{B}_{ij}^{rr'}) \\ &= -\frac{\exp\{-(x_{ij}^{r}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})}{1+\exp\{-(x_{ij}^{r}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}} \operatorname{pr}(X_{ij}^{r}=x_{ij}^{r'},X_{ij}^{r'}=x_{ij}^{r'}\mid \mathcal{B}_{ij}^{rr'}) \\ &+\frac{(x_{ij}^{r}-x_{ij}^{r'})(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})}{1+\exp\{-(x_{ij}^{r}-x_{ij}^{r'})\beta_{j,O\setminus j}^{T}(x_{i,O\setminus j}^{r}-x_{i,O\setminus j}^{r'})\}} \operatorname{pr}(X_{ij}^{r}=x_{ij}^{r'}\mid \mathcal{B}_{ij}^{rr'}). \end{split}$$

Substituting (S8) and (S9) into the last expression yields  $E\{h_{ij}^{rr'}(\beta_{j,O\setminus j}^*) \mid \mathcal{B}_{ij}^{rr'}\} = 0$ . By the law of iterated expectation, we obtain  $E\{h_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\} = 0$ . To show  $E\{\nabla \ell_{i}(\beta^*) = 0\} = 0$  we simply recall from (S6) that

To show  $E\{\nabla \ell_j(\beta^*_{j,O\setminus j})\}=0$ , we simply recall from (S6) that

$$\nabla \ell_j(\beta_{j,O\setminus j}^*) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{2}{R(R-1)} \sum_{r < r'} h_{ij}^{rr'}(\beta_{j,O\setminus j}^*) \right\}$$

and use the fact that  $E\{h^{rr'}_{ij}(\beta^*_{j,O\backslash j})\}=0.$ 

Next, we show that  $\beta_{j,O\setminus j}^*$  is a global minimizer of  $E\{\ell_j(\beta_{j,O\setminus j})\}$ . Each component of  $h_{ij}^{rr'}(\beta_{j,O\setminus j})$  is dominated by the corresponding one of  $|(x_{ij}^r - x_{ij}^{r'})(x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'})|$ , and the latter is integrable by Proposition S1 in § S4·1. Applying the dominated convergence theorem, we can interchange the order of integration and differentiation. Thus, we obtain  $\nabla E\{\ell_j(\beta_{j,O\setminus j})\}|_{\beta_{j,O\setminus j}=\beta_{j,O\setminus j}^*} = E\{\nabla \ell_j(\beta_{j,O\setminus j}^*)\} = 0$ . Hence  $\beta_{j,O\setminus j}^*$  is a stationary point of  $E\{\ell_j(\beta_{j,O\setminus j})\}$ . Following a similar argument, we have

$$\nabla^2 E\{\ell_j(\beta_{j,O\setminus j})\} = E\left[\frac{R_{ij}^{rr'}(\beta_{j,O\setminus j})(X_{ij}^r - X_{ij}^{r'})^2(X_{i,O\setminus j}^r - X_{i,O\setminus j}^{r'})^{\otimes 2}}{\{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j})\}^2}\right] \succeq 0.$$

Therefore  $E\{\ell_j(\beta_{j,O\setminus j})\}$  is convex. It follows that  $\beta_{j,O\setminus j}^*$  is a global minimizer of  $E\{\ell_j(\beta_{j,O\setminus j})\}$ .

#### S4. PROOF OF THEOREMS 1 AND 2

S4-1. *Some technical lemmas* 

Assumption 3 allows us to control the tail behaviour of the random variables. Given Assumption 3, we have the following proposition on the tail probability of the random variables.

**PROPOSITION S1.** Under Assumption 3, for any t > 0 and  $j \in O$  we have

 $\operatorname{pr}\left(|X_j| \ge t\right) \le c_1 \exp(-t),$ 

where  $c_1 = 2 \exp(\kappa_m + \kappa_h/2)$ . Moreover, for two replicates of  $X_j$   $(j \in O)$ , which we denote by  $X_j^r$  and  $X_j^{r'}$ , we have

$$\operatorname{pr}\left(|X_j^r - X_j^{r'}| \ge t\right) \le 2c_1 \exp(-t/2).$$

Also,

$$E(X_j^4) = \int \operatorname{pr}(X_j^4 \ge t) \, \mathrm{d}t \le \int c_1 \exp(-t^{1/4}) \, \mathrm{d}t = 24 \, c_1$$

and

$$E(X_j^8) = \int \operatorname{pr}(X_j^8 \ge t) \, \mathrm{d}t \le \int c_1 \exp(-t^{1/8}) \, \mathrm{d}t = 40\,320\,c_1$$

The proof of Proposition S1 involves the standard Chernoff bounding technique (see Proposition 3 in Yang et al., 2015). We now present a collection of lemmas that will be used to prove Theorems 1 and 2. The proofs of Lemmas S1–S5 are provided in  $\S$  S6.

Recall from (S1) that the gradient of the loss function can be written as the average of independent and identically distributed random variables  $\nabla \ell_j(\beta_{j,O\setminus j}^*) = \sum_{i=1}^n U_i^j(\beta_{j,O\setminus j}^*)/n$ , and that

$$U_{ik}^{j}(\beta_{j,O\setminus j}^{*}) = \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} h_{ijk}^{rr'}(\beta_{j,O\setminus j}^{*}) \in \mathbb{R}.$$
(S10)

The following lemma shows that  $U_i^j(\beta_{i,O\setminus j}^*)$  is a random variable with general exponential tail.

LEMMA S1. For  $j \in O$  and  $k \in O \setminus j$ , we have

$$\operatorname{pr}\left\{U_{ik}^{j}(\beta_{j,O\setminus j}^{*}) \ge t\right\} \leqslant 4c_{1}R^{2}\exp(-0.5t^{1/2}),$$

where  $c_1$  is the constant in Assumption 3.

The following lemma establishes an upper bound for the gradient of the loss function  $\|\nabla \ell_j(\beta_{i,O\setminus j}^*)\|_{\infty}$  with high probability.

LEMMA S2. Under Assumption 3, let

$$\lambda = \frac{(\xi+1)K}{\xi-1} \left(\frac{\log^5 p}{n}\right)^{1/2},$$

where K > 0 and  $\xi > 1$ . Then, for  $j \in O$ ,

$$\left\| \nabla \ell_j(\beta_{j,O\setminus j}^*) \right\|_{\infty} \leq \frac{(\xi-1)\lambda}{\xi+1}$$

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

5

85

90

95

100

Now, let  $D(\hat{\beta}_{j,O\setminus j}, \beta_{j,O\setminus j}) = (\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j})^{\mathrm{T}} \{\nabla \ell_j(\hat{\beta}_{j,O\setminus j}) - \nabla \ell_j(\beta_{j,O\setminus j})\}$  be the symmetrized Bregman divergence of  $\ell_j(\beta_{j,O\setminus j})$ . The next lemma allows us to control the variability of the Bregman divergence and the Hessian matrix in a neighbourhood of  $\beta_{j,O\setminus j}$ .

LEMMA S3. Let  $\Delta \in \mathbb{R}^{p-1}$  and

$$b = \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left| (x_{ij}^r - x_{ij}^{r'}) \Delta^{\mathrm{T}} (x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'}) \right|.$$

Then, for  $j \in O$ ,

$$\exp(-b)\nabla^2 \ell_j(\beta_{j,O\setminus j}) \leqslant \nabla^2 \ell_j(\beta_{j,O\setminus j} + \Delta) \leqslant \exp(b)\nabla^2 \ell_j(\beta_{j,O\setminus j})$$

and

$$\exp(-b)\Delta^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j})\Delta \leqslant D(\beta_{j,O\setminus j}+\Delta,\beta_{j,O\setminus j}) \leqslant \exp(b)\Delta^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j})\Delta.$$

We omit the proof, but the result follows from arguments similar to those in Ning et al. (2016). Given Lemma S3, we have the following inequality.

LEMMA S4. Let  $\hat{\Delta}_j = \hat{\beta}_{j,O\setminus j} - \beta^*_{j,O\setminus j}$ . Then, for  $j \in O$ ,  $D(\hat{\beta}_{j,O\setminus j}, \beta^*_{j,O\setminus j}) + \{\lambda - \|\nabla \ell_j(\beta^*_{j,O\setminus j})\|_{\infty}\} \|\hat{\Delta}_{j,\mathcal{S}_j^c}\|_1 \leq \{\lambda + \|\nabla \ell_j(\beta^*_{j,O\setminus j})\|_{\infty}\} \|\hat{\Delta}_{j,\mathcal{S}_j}\|_1$ . Moreover, for any  $\xi > 1$ ,  $\|\hat{\Delta}_{j,\mathcal{S}_i^c}\|_1 \leq \xi \|\hat{\Delta}_{j,\mathcal{S}_j}\|_1$ , provided that  $\|\nabla \ell_j(\beta^*_{j,O\setminus j})\|_{\infty} \leq \lambda(\xi - \xi)$ .

 $1)/(\xi+1).$ 

Recall from §2.4 that  $H^j = E\{\nabla^2 \ell_j(\beta^*_{j,O\setminus j})\}$ . We next present a lemma on the deviation between the empirical Hessian matrix and the population Hessian matrix.

LEMMA S5. Under Assumption 3, for  $j \in O$  there exists a constant  $K_1 > 0$  such that

$$\left\|\nabla^2 \ell_j(\beta_{j,O\setminus j}^*) - H^j\right\|_{\infty} \leqslant K_1 \left(\log^9 p/n\right)^{1/2}$$

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

# S4.2. Proof of Theorem 1

To obtain an upper bound for the estimation error, we state some conditions on the minimal eigenvalues of the Hessian matrix of the loss function. Given the constant  $\xi > 1$ , for each  $j \in O$ we define the compatibility factor, restricted eigenvalue and weak cone invertibility factor as

$$\kappa^{2}\left\{\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*}), s_{j}\right\} = \min_{v}\left\{\frac{s_{j}v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v}{\|v_{S_{j}}\|_{1}^{2}} : v \in \mathbb{R}^{p-1}, v \neq 0, \|v_{S_{j}^{c}}\|_{1} \leqslant \xi \|v_{S_{j}}\|_{1}\right\},$$
(S11)

$$\operatorname{RE}\left\{\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*}), s_{j}\right\} = \min_{v} \left\{\frac{v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v}{\|v\|_{2}^{2}} : v \in \mathbb{R}^{p-1}, v \neq 0, \|v_{S_{j}^{c}}\|_{1} \leqslant \xi \|v_{S_{j}}\|_{1}\right\}$$
(S12)

and, for  $q \ge 1$ ,

$$\rho_q \left\{ \nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j \right\} = \min_v \left\{ \frac{s_j^{1/q} v^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*) v}{\|v_{S_j}\|_1 \|v\|_q} : v \in \mathbb{R}^{p-1}, v \neq 0, \|v_{S_j^c}\|_1 \leqslant \xi \|v_{S_j}\|_1 \right\}.$$
(S13)

These conditions have commonly been used to establish upper bounds for estimation error in the context of  $\ell_1$ -penalized regression (Bickel et al., 2009; van de Geer & Bühlmann, 2009; Ye & Zhang, 2010). Define  $\kappa_{\min}^2 = \min_{j \in O} \kappa^2 \{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}$ , RE<sub>min</sub> =  $\min_{j \in O} \operatorname{RE}\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}$  and  $\rho_{q,\min} = \min_{j \in O} \rho_q\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}$ . We now give the proof of Theorem 1.

Proof of Theorem 1. The proof involves obtaining an upper bound on the gradient of the loss function. By Lemma S2 we have that  $\|\nabla \ell_j(\beta_{j,O\setminus j}^*)\|_{\infty} \leq \lambda(\xi-1)/(\xi+1)$  with probability at least  $1-p^{-1}$ . Throughout the proof, we conditioned on the event  $\|\nabla \ell_j(\beta_{j,O\setminus j}^*)\|_{\infty} \leq \lambda(\xi-1)/(\xi+1)$ .

Let 
$$\Delta_j = \beta_{j,O\setminus j} - \beta_{j,O\setminus j}^*$$
,  $a = \Delta_j / \|\Delta_j\|_1$  and  $k = \|\Delta_j\|_1$ . Recall from Lemma S3 that  
 $b = \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left| (x_{ij}^r - x_{ij}^{r'}) \hat{\Delta}_j^{\mathrm{T}} (x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'}) \right|.$ 

Consider the quantity  $D(\hat{\beta}_{j,O\setminus j}, \beta^*_{j,O\setminus j})$ . By Lemma S3,

$$D(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*) \ge \exp(-b)k^2 a^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*)a \ge k^2 \exp(-Mk)a^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*)a,$$
(S14)

since

$$b = \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left\| (x_{ij}^r - x_{ij}^{r'}) \hat{\Delta}_j^{\mathrm{T}} (x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'}) \right\|$$

$$\leq \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left\| (x_{ij}^r - x_{ij}^{r'}) (x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'}) \right\|_{\infty} \| \hat{\Delta}_j \|_1$$

$$\leq Mk,$$

$$(14)$$

where the last inequality holds conditioning on the event

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

By the definition of the compatibility factor, (S11), we have

$$a^{\mathrm{T}} \nabla^2 \ell_j(\beta^*_{j,O\setminus j}) a \ge \kappa^2 \{ \nabla^2 \ell_j(\beta^*_{j,O\setminus j}), s_j \} \|a_{\mathcal{S}_j}\|_1^2 / s_j.$$

Substituting this into (S14) gives

$$D\left(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*\right) \geqslant k^2 \exp(-Mk) \kappa^2 \{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\} \|a_{\mathcal{S}_j}\|_1^2 / s_j.$$
(S15)

Next, we derive an upper bound for  $D(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*)$ . By Lemma S4, conditioning on the event  $\|\nabla \ell_j(\beta_{j,O\setminus j}^*)\|_{\infty} \leq \lambda(\xi-1)/(\xi+1)$ , we have  $\|\hat{\Delta}_{j,S_j^c}\| \leq \xi \|\hat{\Delta}_{j,S_j}\|$ . Hence,

$$D(\beta_{j,O\setminus j}^{*} + ka, \beta_{j,O\setminus j}^{*}) \leq \{\lambda + \|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty}\} \|\hat{\Delta}_{j,S_{j}}\|_{1} - \{\lambda - \|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty}\} \|\hat{\Delta}_{j,S_{j}^{c}}\|_{1} \\ = k\{\lambda + \|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty}\} \|a_{S_{j}}\|_{1} - k\{\lambda - \|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty}\} \|a_{S_{j}^{c}}\|_{1} \\ \leq \frac{2k\xi\lambda}{\xi+1} \|a_{S_{j}}\|_{1} - \frac{2k\lambda}{\xi+1} \|a_{S_{j}^{c}}\|_{1} + \frac{2k\lambda}{\xi+1} \|a_{S_{j}}\|_{1} - \frac{2k\lambda}{\xi+1} \|a_{S_{j}}\|_{1}$$

$$= 2k\lambda \|a_{S_{j}}\|_{1} - \frac{2k\lambda}{\xi+1} \\ \leq \frac{k\lambda(\xi+1)\|a_{S_{j}}\|_{1}^{2}}{2}, \qquad (S16)$$

130

where the second equality follows from the fact that  $||a||_1 = 1$  by construction, and the last inequality is obtained by using the fact that  $a^2 + b^2 \ge 2ab$  for  $a, b \in \mathbb{R}$ .

<sup>155</sup> Combining (S15) and (S16), we have

$$k \exp(-Mk) \leqslant \frac{\xi + 1}{2\kappa^2 \{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}} \lambda s_j.$$

Let  $\tau = M(\xi + 1)\lambda s_j / [2\kappa^2 \{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}]$ . Then we have  $Mk \exp(-Mk) \leq \tau$ . Since  $\eta$  is the smallest solution of  $z \exp(-z) = \tau$  by definition, and  $z \exp(-z) - \tau$  is an increasing function of z for  $z \leq 1$ , this implies that  $Mk \leq \eta$ . Therefore

$$\|\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j}^*\|_1 = \|\hat{\Delta}_j\|_1 = k \leqslant \frac{\eta}{M} = \frac{\tau \exp(\eta)}{M} = \frac{(\xi+1)\exp(\eta)}{2\kappa^2 \{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}} \lambda s_j,$$

which implies  $\|\hat{\beta}_{j,O\setminus j} - \beta^*_{j,O\setminus j}\|_1 \leq (\xi+1)\exp(\eta)s_{\max}\lambda/(2\kappa_{\min}^2)$ , where  $\kappa_{\min}^2 = \min_{j\in O} \kappa^2 \{\nabla^2 \ell_j(\beta^*_{j,O\setminus j}), s_j\}.$ 

To prove  $\|\hat{\beta}_{j,O\setminus j} - \hat{\beta}^*_{j,O\setminus j}\|_2 \leq C'(s_{\max})^{1/2}\lambda/\operatorname{RE}_{\min}$ , we recall that by the definition of the restricted eigenvalue, (S12),  $a^{\mathrm{T}}\nabla^2\ell_j(\beta^*_{j,O\setminus j})a \geq \operatorname{RE}\{\nabla^2\ell_j(\beta^*_{j,O\setminus j}), s_j\}\|a\|_2^2$ . Thus, by (S14),

$$D\left(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*\right) \geqslant k^2 \exp(-Mk) \operatorname{RE}\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\} \|a\|_2^2.$$
(S17)

Similar to (S16), by an application of Lemma S4 it can be shown that

$$D(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*) \leqslant \frac{2k\xi\lambda}{1+\xi} \|a_{\mathcal{S}_j}\|_1 \leqslant \frac{2k\xi\lambda}{1+\xi} s_j^{1/2} \|a_{\mathcal{S}_j}\|_2 \leqslant \frac{2k\xi\lambda}{1+\xi} s_j^{1/2} \|a\|_2.$$
(S18)

Upon combining (S17) and (S18), we have

$$\|\hat{\beta}_{j,O\backslash j} - \beta^*_{j,O\backslash j}\|_2 = k\|a\|_2 \leqslant \frac{2\xi \exp(\eta)}{(\xi+1)\operatorname{RE}\{\nabla^2 \ell_j(\beta^*_{j,O\backslash j}), s_j\}} s_j^{1/2}\lambda \leqslant \frac{2\xi \exp(\eta)}{(\xi+1)\operatorname{RE}_{\min}} s_{\max}^{1/2}\lambda,$$

165

where the first inequality holds because  $Mk \leq \eta$  and  $\operatorname{RE}_{\min} = \min_{j \in O} \operatorname{RE}\{\nabla^2 \ell_j(\beta^*_{j,O\setminus j}), s_j\}$ .

To prove  $\|\hat{\beta}_{j,O\setminus j} - \beta^*_{j,O\setminus j}\|_q \leq C'(s_{\max})^{1/q}\lambda/\rho_{q,\min}$ , recall that by the definition of the weak cone invertibility factor, (S13), we have  $a^{\mathrm{T}}\nabla^2\ell_j(\beta^*_{j,O\setminus j})a \geq \rho_q\{\nabla^2\ell_j(\beta^*_{j,O\setminus j}), s_j\}\|a_{\mathcal{S}_j}\|_1\|a\|_q/s_j^{1/q}$ . Hence, by (S14),

$$D(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*) \ge k^2 \exp(-Mk)\rho_q\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\} \|a_{\mathcal{S}_j}\|_1 \|a\|_q / s_j^{1/q}.$$
 (S19)

Moreover, from (S18) we have

$$D\left(\beta_{j,O\setminus j}^* + ka, \beta_{j,O\setminus j}^*\right) \leqslant \frac{2k\xi\lambda}{1+\xi} \|a_{\mathcal{S}_j}\|_1.$$
(S20)

170 Combining (S19) and (S20) yields

$$\left\|\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j}^*\right\|_q = k\|a\|_q \leqslant \frac{2\xi \exp(\eta)}{(\xi+1)\rho_q\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j\}} s_j^{1/q} \lambda \leqslant \frac{2\xi \exp(\eta)}{(\xi+1)\rho_{q,\min}} s_{\max}^{1/q} \lambda,$$

where the first inequality follows from the facts that  $Mk \leq \eta$  and  $\rho_{q,\min} = \min_{j \in O} \rho_q \{ \nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j \}.$ 

# S4-3. Proof of Theorem 2

Let  $\mathcal{A} = (v \in \mathbb{R}^{p-1} : \|v_{S_j^c}\|_1 \leq \xi \|v_{S_j}\|_1)$ . We first show that for any  $v \in \mathcal{A}$ , the quantity  $v^{\mathrm{T}} E\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)\}v$  can be bounded below by a positive constant. Using the fact that the deviation between the empirical and population Hessian matrices is arbitrarily small when n is sufficiently large (see Lemma S5), we show that  $v^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*)v$  is also bounded below by a constant. Therefore the compatibility factor (S11), the restricted eigenvalue (S12), and the weak cone invertibility factor (S13) can all be bounded below by positive constants.

Lower bound for  $v^{\mathrm{T}}E\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)\}v$ . We show this via a truncation argument. For 180  $\delta, S_1, S_2 > 0$  and  $j \in O$ , we define the events

$$\begin{split} F_{ij}^{rr'} &= \left( |x_{ij}^r| \leqslant \delta \right) \cap \left( |x_{ij}^{r'}| \leqslant \delta \right), \\ G_{ij}^r &= \left\{ \left| (\beta_{j,O\setminus j}^*)^{\mathrm{T}} x_{i,O\setminus j}^r \right| \leqslant S_1 \right\}, \qquad G_{ij}^{r'} &= \left\{ \left| (\beta_{j,O\setminus j}^*)^{\mathrm{T}} x_{i,O\setminus j}^{r'} \right| \leqslant S_1 \right\}, \\ H_{ij}^r &= \left\{ \left| (\beta_{j,-j}^*)^{\mathrm{T}} x_{i,-j}^r \right| \leqslant S_2 \right\}, \qquad H_{ij}^{r'} &= \left\{ \left| (\beta_{j,-j}^*)^{\mathrm{T}} x_{i,-j}^{r'} \right| \leqslant S_2 \right\}. \end{split}$$

Since  $\exp(-z)/\{1 + \exp(-z)\}^2$  is a decreasing function of z for  $z \ge 0$ , the quantity  $R_{ij}^{rr'}(\beta_{j,O\setminus j}^*)/\{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\}^2$  can be bounded below by  $C_1 = \exp(-4S_1\delta)/\{1 + \exp(-4S_1\delta)\}^2$ . Recall from (S4) and (S7) the definitions of  $T_{ij}^{rr'}$  and  $\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)$ . We have

$$\begin{split} \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) \\ \geqslant \frac{1}{n} \sum_{i=1}^{n} \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} T_{ij}^{rr'} I(F_{ij}^{rr'}) I(G_{ij}^{r}) I(G_{ij}^{r'}) \\ \geqslant \frac{1}{n} \sum_{i=1}^{n} \frac{2}{R(R-1)} \sum_{1 \leqslant r < r' \leqslant R} C_{1} (x_{ij}^{r} - x_{ij}^{r'})^{2} (x_{i,O\setminus j}^{r} - x_{i,O\setminus j}^{r'})^{\otimes 2} I(F_{ij}^{rr'}) I(G_{ij}^{r}) I(G_{ij}^{r'}) \\ = \frac{1}{n} \sum_{i=1}^{n} W_{ij}. \end{split}$$

Let

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

and let  $\Theta = \Sigma^{-1}$ . Assume that  $X_{iH} \sim N(0, \Sigma_{H,H})$  and that  $X_{iO}^r \mid X_{iH}$  are independent and identically distributed from  $N(\Sigma_{O,H}\Sigma_{H,H}^{-1}X_{iH}, \Sigma_{O,O} - \Sigma_{O,H}\Sigma_{H,H}^{-1}\Sigma_{H,O})$  for  $r = 1, \ldots, R$ . <sup>190</sup> Then

$$\boldsymbol{\rho}(x_{ij}^r \mid x_{i,-j}^r, H_{ij}^r) = \frac{\boldsymbol{\rho}(x_i^r \mid H_{ij}^r)}{\int \boldsymbol{\rho}(x_i^r \mid H_{ij}^r) \, \mathrm{d}x_{ij}}$$
$$= \frac{\boldsymbol{\rho}(x_i^r)}{\boldsymbol{\rho}(H_{ij}^r) \int \boldsymbol{\rho}(x_i^r) / \boldsymbol{\rho}(H_{ij}^r) \, \mathrm{d}x_{ij}}$$
$$= \boldsymbol{\rho}(x_{ij}^r \mid x_{i,-j}^r),$$

where we have used the facts that  $p(x_i^r \mid H_{ij}^r) = p(x_i^r)/p(H_{ij}^r)$  and  $p(H_{ij}^r)$  is a constant. Recall from Example 1 that the conditional density of  $X_{ij}^r$  given  $X_{i,-j}^r$  is

$$\boldsymbol{\rho}(x_{ij}^r \mid x_{i,-j}^r) = (\Theta_{jj}/2\pi)^{1/2} \exp\left\{-\frac{\Theta_{jj}}{2}(x_{ij}^r)^2 - x_{ij}^r \Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^r - \frac{1}{2\Theta_{jj}} \left(\Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^r\right)^2\right\},$$

where  $\beta_{j,-j}^* = -\Theta_{j,-j}$ . Since  $X_{ij}^r$  and  $X_{ij}^{r'}$  are conditionally independent and identically normally distributed random variables, we can bound the following conditional expectation from below:

$$\begin{split} E\Big\{ (X_{ij}^r - X_{ij}^{r'})^2 I(F_{ij}^{rr'}) \mid x_{i,-j}^r, x_{i,-j}^{r'}, H_{ij}^r, H_{ij}^{r'} \Big\} \\ &= \Theta_{jj}/2\pi \int_{-\delta}^{\delta} \int_{-\delta}^{\delta} (x_{ij}^r - x_{ij}^{r'})^2 \exp\left\{ -\frac{\Theta_{jj}}{2} (x_{ij}^r)^2 - x_{ij}^r \Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^r - \frac{1}{2\Theta_{jj}} \left(\Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^r\right)^2 \right\} \\ &\qquad \times \exp\left\{ -\frac{\Theta_{jj}}{2} (x_{ij}^{r'})^2 - x_{ij}^{r'} \Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^{r'} - \frac{1}{2\Theta_{jj}} \left(\Theta_{j,-j}^{\mathrm{T}} x_{i,-j}^{r'}\right)^2 \right\} \mathrm{d}x_{ij}^r \mathrm{d}x_{ij}^{r'} \\ &\geqslant \Theta_{jj}/2\pi \int_{-\delta}^{\delta} \int_{-\delta}^{\delta} (x_{ij}^r - x_{ij}^{r'})^2 \\ &\qquad \times \exp\left[ -\frac{\Theta_{jj}}{2} \left\{ (x_{ij}^r)^2 + (x_{ij}^{r'})^2 \right\} - S_2(x_{ij}^r + x_{ij}^{r'}) - \frac{S_2^2}{2\Theta_{jj}} \right] \mathrm{d}x_{ij}^r \mathrm{d}x_{ij}^{r'}, \end{split}$$

where the last inequality follows from the fact that conditioned on the events  $H_{ij}^r$  and  $H_{ij}^{r'}$ ,  $|(\beta_{j,-j}^*)^{\mathrm{T}} x_{i,-j}^r| \leq S_2$ . For notational convenience, we denote the last expression by  $C_2$ .

Therefore, by the law of iterated expectation, we obtain

$$v^{\mathrm{T}} E \{ \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) \} v \ge v^{\mathrm{T}} E \left( \frac{1}{n} \sum_{i=1}^{n} W_{ij} \right) v$$

$$= v^{\mathrm{T}} E(W_{ij}) v$$

$$= v^{\mathrm{T}} E \{ E(W_{ij} \mid x_{i,-j}^{r} x_{i,-j}^{r'}, H_{ij}^{r}, H_{ij}^{r'} \forall r < r') \} v$$

$$\ge C_{1} C_{2} E \left[ \{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}} v \}^{2} I(G_{ij}^{r}) I(G_{ij}^{r'}) \right],$$
(S21)

where we have used the fact that the replicates are conditionally independent and identically distributed. We now establish a lower bound for (S21). By the Cauchy–Schwarz inequality, we have

$$E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v \right\}^{2} \left\{ 1 - I(G_{ij}^{r})I(G_{ij}^{r'}) \right\} \right] \\ \leqslant \left( E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v \right\}^{4} \right] \right)^{1/2} \left[ \operatorname{pr}\left\{ (G_{ij}^{r})^{\mathrm{c}} \cup (G_{ij}^{r'})^{\mathrm{c}} \right\} \right]^{1/2}.$$
(S22)

Recall that  $X_{iO}^r \mid X_{iH}$  are independent and identically distributed from  $N(\Sigma_{O,H}\Sigma_{H,H}^{-1}X_{iH}, \Sigma_{O,O} - \Sigma_{O,H}\Sigma_{H,H}^{-1}\Sigma_{H,O})$ . For notational convenience, we write  $\Sigma' = \Sigma_{O,O} - \Sigma_{O,H}\Sigma_{H,H}^{-1}\Sigma_{H,O}$ . Since  $X_{i,O\setminus j}^r \mid X_{iH} \sim N(\Sigma_{O\setminus j,H}\Sigma_{H,H}^{-1}X_{iH}, \Sigma'_{O\setminus j,O\setminus j})$ , we have that  $(X_{i,O\setminus j}^r - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v \mid X_{iH} \sim N(0, 2v^{\mathrm{T}}\Sigma'_{O\setminus j,O\setminus j}v)$ . Therefore, the kurtosis of a normal distribution is

$$E\left[\left\{ (X_{i,O\setminus j}^r - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v\right\}^4 \right] = E\left(E\left[\left\{ (X_{i,O\setminus j}^r - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v\right\}^4 \mid X_{iH}\right]\right)$$
$$= 3(2v^{\mathrm{T}}\Sigma_{O\setminus j,O\setminus j}'v)^2$$
$$\leqslant 12||v||_2^4 \Lambda_{\max}^2(\Sigma),$$

where  $\Lambda_{\max}(\Sigma)$  is the largest eigenvalue of  $\Sigma$ . The last inequality is obtained from the fact that  $\Lambda_{\max}(\Sigma') \leq \Lambda_{\max}(\Sigma)$ .

Recall that  $X_{i,O\setminus j}^r \sim N(0, \Sigma_{O\setminus j,O\setminus j})$ . Therefore  $(\beta_{j,O\setminus j}^*)^{\mathrm{T}} X_{i,O\setminus j}^r \sim N(0, \sigma_1^2)$ , where  $\sigma_1^2 = \beta_{j,O\setminus j}^* \sum_{O\setminus j,O\setminus j} \beta_{j,O\setminus j}^*$ . By the Gaussian tail inequality in Lemma S13, we have

$$\Pr\{(G_{ij}^r)^{c} \cup (G_{ij}^{r'})^{c}\} \leq 2\Pr\{(G_{ij}^r)^{c}\} = 2\Pr\{\left|(\beta_{j,O\setminus j}^*)^{T}X_{i,O\setminus j}^r\right| > S_1\} \leq \frac{4\sigma_1}{S_1}\exp(-S_1^2/2\sigma_1^2).$$

We write the last expression as  $C_3(S_1)$ , indicating its dependence on  $S_1$ . Therefore, by (S22) and picking a sufficiently large  $S_1$  such that  $(12)^{1/2} \Lambda_{\max}(\Sigma) \{C_3(S_1)\}^{1/2} = \Lambda_{\min}(\Sigma)$ , we obtain

$$E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v \right\}^{2} \left\{ 1 - I(G_{ij}^{r})I(G_{ij}^{r'}) \right\} \right]$$
  
$$\leq (12)^{1/2} \|v\|_{2}^{2} \Lambda_{\max}(\Sigma) \{C_{3}(S_{1})\}^{1/2} = \|v\|_{2}^{2} \Lambda_{\min}(\Sigma).$$
(S23)

In addition, by Hölder's inequality, we have

$$E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v\right\}^{2}\right] = E\left(E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v\right\}^{2} \mid X_{iH}\right]\right)$$
  
$$= 2v^{\mathrm{T}}\Sigma_{O\setminus j,O\setminus j}^{\prime}v$$
  
$$\geq 2\|v\|_{2}^{2}\Lambda_{\min}(\Sigma).$$
(S24)

Combining (S23) and (S24), we obtain

$$2\|v\|_{2}^{2}\Lambda_{\min}(\Sigma) \leqslant \|v\|_{2}^{2}\Lambda_{\min}(\Sigma) + E\left[\left\{(X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v\right\}^{2}I(G_{ij}^{r})I(G_{ij}^{r'})\right],$$

implying

$$E\left[\left\{ (X_{i,O\setminus j}^{r} - X_{i,O\setminus j}^{r'})^{\mathrm{T}}v \right\}^{2} I(G_{ij}^{r}) I(G_{ij}^{r'}) \right] \ge \|v\|_{2}^{2} \Lambda_{\min}(\Sigma).$$
(S25)

By (S21) and (S25), we conclude that

$$v^{\mathrm{T}} E\{\nabla^2 \ell_j(\beta_{j,O\backslash j}^*)\} v \ge C_1 C_2 \|v\|_2^2 \Lambda_{\min}(\Sigma).$$
(S26)

Lower bound for  $v^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*)v$ . Gaussian random variables satisfy Assumption 3 (see Yang et al., 2015). Also, recall that  $H^j = E\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)\}$ . Therefore, by Lemma S5,

$$\left\|\nabla^2 \ell_j(\beta_{j,O\setminus j}^*) - H^j\right\|_{\infty} \leqslant K_1(\log^9 p/n)^{1/2}$$

with probability  $1 - p^{-1}$ . Let  $\Gamma = H^j - \nabla^2 \ell_j(\beta^*_{j,O\setminus j})$ . By Hölder's inequality,

$$v^{\mathrm{T}}H^{j}v - v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v \leqslant v^{\mathrm{T}}\Gamma v \leqslant \|v\|_{1}^{2}\|\Gamma\|_{\infty} \leqslant \|v\|_{1}^{2}K_{1}(\log^{9}p/n)^{1/2}.$$
 (S27)

Note that  $||v_{\mathcal{S}_j}||_1 \leq s_j^{1/2} ||v_{\mathcal{S}_j}||_2 \leq s_j^{1/2} ||v||_2$  and that for any  $v \in \mathcal{A}$ ,  $||v_{\mathcal{S}_j^c}||_1 \leq \xi ||v_{\mathcal{S}_j}||_1$ . By combining (S26) and (S27) and using the above facts, we obtain

$$C_{1}C_{2}\|v\|_{2}^{2}\Lambda_{\min}(\Sigma) \leq \|v\|_{1}^{2}K_{1}(\log^{9}p/n)^{1/2} + v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v$$
  
$$\leq (1+\xi)^{2}\|v_{\mathcal{S}_{j}}\|_{1}^{2}K_{1}(\log^{9}p/n)^{1/2} + v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v$$
  
$$\leq (1+\xi)^{2}\|v\|_{2}^{2}K_{1}s_{j}(\log^{9}p/n)^{1/2} + v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v.$$

By the assumption  $\lim_{n\to\infty} s_{\max} (\log^9 p/n)^{1/2} = 0$ , for sufficiently large n we have

$$(1+\xi)^2 K_1 s_j (\log^9 p/n)^{1/2} \leqslant \frac{1}{2} C_1 C_2 \Lambda_{\min}(\Sigma).$$

215

Hence

$$\frac{v^{\mathrm{T}} \nabla^2 \ell_j(\beta^*_{j,O\setminus j})v}{\|v\|_2^2} \geqslant \frac{1}{2} C_1 C_2 \Lambda_{\min}(\Sigma).$$
(S28)

*Compatibility factor.* By the definition of the compatibility factor, (S11), together with (S28), we have

$$\kappa^{2} \left\{ \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}), s_{j} \right\} = \min_{v \in \mathcal{A}} \frac{s_{j} v^{\mathrm{T}} \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) v}{\|v_{\mathcal{S}_{j}}\|_{1}^{2}}$$
$$\geqslant \min_{v \in \mathcal{A}} \frac{v^{\mathrm{T}} \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) v}{\|v\|_{2}^{2}} \geqslant \frac{1}{2} C_{1} C_{2} \Lambda_{\min}(\Sigma),$$

where we have used the fact that  $\|v_{\mathcal{S}_j}\|_1 \leqslant s_j^{1/2} \|v_{\mathcal{S}_j}\|_2 \leqslant s_j^{1/2} \|v\|_2$ .

*Restricted eigenvalue.* By the definition of the restricted eigenvalue, (S12), together with (S28), we have

$$\operatorname{RE}\left\{\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*}), s_{j}\right\} = \min_{v \in \mathcal{A}} \frac{v^{\mathrm{T}}\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})v}{\|v\|_{2}^{2}} \geqslant \frac{1}{2}C_{1}C_{2}\Lambda_{\min}(\Sigma).$$

*Weak cone invertibility factor.* By the definition of the weak cone invertibility factor, (S13), together with (S28), we have

235

$$\rho_q \left\{ \nabla^2 \ell_j(\beta_{j,O\setminus j}^*), s_j \right\} = \min_{v \in \mathcal{A}} \frac{s_j^{1/q} v^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*) v}{\|v_{\mathcal{S}_j}\|_1 \|v\|_q}$$
$$\geqslant \min_{v \in \mathcal{A}} \frac{s_j v^{\mathrm{T}} \nabla^2 \ell_j(\beta_{j,O\setminus j}^*) v}{\|v_{\mathcal{S}_j}\|_1^2} \geqslant \frac{1}{2} C_1 C_2 \Lambda_{\min}(\Sigma),$$

where the first inequality follows from Hölder's inequality, i.e.,  $\|v_{S_j}\|_1 \leq s_j^{1-1/q} \|v\|_q$ .

## S5. PROOF OF THEOREM 3

We start by presenting some lemmas that will be used to prove Theorem 3. The proofs of these lemmas are provided in § S7. Recall from (9) that the pairwise decorrelated score function for  $\beta_{jk}$  is defined as

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

The estimated pairwise decorrelated score function defined in (11) is

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

The following lemma establishes the asymptotic normality of  $S_{jk}(\beta_{j\vee k})$ .

LEMMA S6. Under Assumptions 3–5, for  $j, k \in O$  we have that  $n^{1/2}S_{jk}(\beta_{j\vee k}^*)$  converges in distribution to  $N(0, \sigma_{jk}^2)$ , where

$$\sigma_{jk}^2 = \Sigma_{jk,jk}^{jk} - 2\Sigma_{jk,j\backslash k}^{jk} w_{jk}^* - 2\Sigma_{jk,k\backslash j}^{jk} w_{kj}^* + (w_{jk}^*)^{\mathrm{T}} \Sigma_{j\backslash k,j\backslash k}^{jk} w_{jk}^* + (w_{kj}^*)^{\mathrm{T}} \Sigma_{k\backslash j,k\backslash j}^{jk} w_{kj}^*.$$

Lemma S6 shows that the pairwise decorrelated score function converges to a univariate Gaussian distribution. To derive the asymptotic distribution of  $\hat{S}_{jk}$ , we show that  $n^{1/2}\{\hat{S}_{jk} - S_{jk}(\beta_{j\vee k}^*)\} = o_{\mathbb{P}}(1)$ . We then use Lemma S6 to establish the asymptotic normality of  $\hat{S}_{jk}$ . To this end, we need some additional assumptions on the scaling of n and p and on the magnitude of the regularization parameters  $\lambda$  and  $\lambda_w$  in (7) and (10), respectively, as stated in Assumption S1. 250 We also need the following technical lemmas.

LEMMA S7. Under Assumptions 3 and 4, for  $j \in O$  we have

$$\left\|\nabla^2 \ell_j(\hat{\beta}_{j,O\setminus j}) - H^j\right\|_{\infty} = \mathcal{O}_{\mathbb{P}}\left\{\frac{Ms_{\max}\lambda}{\kappa_{\min}^2} + \left(\frac{\log^9 p}{n}\right)^{1/2}\right\},\$$

where M is a constant from Theorem 1.

LEMMA S8. Let  $w_0 = \max_{j,k \in O} \|w_{jk}^*\|_1$ , and let

$$\lambda_w \ge C \left[ w_0 \left\{ \frac{M s_{\max} \lambda}{\kappa_{\min}^2} + \left( \frac{\log^9 p}{n} \right)^{1/2} \right\} \right],$$

where M is as defined in Theorem 1 and C > 0 is a sufficiently large constant. Under Assumptions 3 and 4, for  $j, k \in O$ ,

$$\left\|\nabla_{k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j}) - (w_{jk}^{*})^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j})\right\|_{\infty} \leq \lambda_{w}$$

with probability converging to 1.

LEMMA S9. Let  $s'_{jk} = \|w_{jk}^*\|_0$  and let  $s'_{\max} = \max_{j,k\in O} s'_{jk}$ . Under Assumptions 3, 4 and S1,  $\|\hat{w}_{jk} - w_{jk}^*\|_1 = \mathcal{O}_{\mathbb{P}}(s'_{\max}\lambda_w)$  holds for all  $j,k\in O$ .

With Lemmas S6–S9, we establish the asymptotic normality of the estimated pairwise decorrelated score function (11).

LEMMA S10. Under Assumptions 3–5 and S1 and under the null hypothesis  $H_0: \beta_{jk}^* = \beta_{kj}^* = 0$  for  $j, k \in O$ , we have that  $n^{1/2}\hat{S}_{jk} = n^{1/2}S_{jk}(\beta_{j\vee k}^*) + o_{\mathbb{P}}(1)$  converges in distribution to  $N(0, \sigma_{jk}^2)$ , where

$$\sigma_{jk}^{2} = \Sigma_{jk,jk}^{jk} - 2\Sigma_{jk,j\backslash k}^{jk} w_{jk}^{*} - 2\Sigma_{jk,k\backslash j}^{jk} w_{kj}^{*} + (w_{jk}^{*})^{\mathrm{T}} \Sigma_{j\backslash k,j\backslash k}^{jk} w_{jk}^{*} + (w_{kj}^{*})^{\mathrm{T}} \Sigma_{k\backslash j,k\backslash j}^{jk} w_{kj}^{*}.$$
(S29)

However,  $\sigma_{jk}^2$  depends on the unknown quantity  $\Sigma^{jk}$ . Recall from (16) that we estimate  $\Sigma^{jk}$  as

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

and that we write  $\hat{\Sigma}^{jk}$  to indicate  $\hat{\Sigma}^{jk}(0, \hat{\beta}_{j,O\setminus\{j,k\}}, \hat{\beta}_{k,O\setminus\{j,k\}})$ . The next lemma asserts that  $\hat{\Sigma}^{jk}$  is a consistent estimator of  $\Sigma^{jk}$ .

LEMMA S11. Let  $\hat{\beta}'_{j\vee k} = (0, \hat{\beta}^{\mathrm{T}}_{j,O\setminus\{j,k\}}, \hat{\beta}^{\mathrm{T}}_{k,O\setminus\{j,k\}})^{\mathrm{T}}$ . Under Assumptions 3, 4 and S1 and the null hypothesis  $H_0: \beta^*_{jk} = \beta^*_{kj} = 0$ ,

$$\hat{\Sigma}^{jk}(\hat{\beta}'_{j\vee k}) = \frac{1}{n} \sum_{i=1}^{n} \left\{ g_i^{jk}(\hat{\beta}'_{j\vee k}) \right\}^{\otimes 2}$$

is a consistent estimator of  $\Sigma^{jk} = E[\{g_i^{jk}(\beta_{j\vee k}^*)\}^{\otimes 2}]$  as defined in (15). In particular,

$$\left\|\hat{\Sigma}^{jk}(\hat{\beta}'_{j\vee k}) - \Sigma^{jk}\right\|_{\infty} = \mathcal{O}_{\mathbb{P}}\left\{\frac{s_{\max}}{\kappa_{\min}^2}\lambda\log^6 p + \left(\frac{\log^9 p}{n}\right)^{1/2}\right\}.$$

With Lemma S11, we establish that  $\hat{\sigma}_{jk}^2$  in (S29) is a consistent estimator of  $\sigma_{jk}^2$ .

LEMMA S12. Let

$$\hat{\sigma}_{jk}^2 = \hat{\Sigma}_{jk,jk}^{jk} - 2\hat{\Sigma}_{jk,j\backslash k}^{jk}\hat{w}_{jk} - 2\hat{\Sigma}_{jk,k\backslash j}^{jk}\hat{w}_{kj} + \hat{w}_{jk}^{\mathrm{T}}\hat{\Sigma}_{j\backslash k,j\backslash k}^{jk}\hat{w}_{jk} + \hat{w}_{kj}^{\mathrm{T}}\hat{\Sigma}_{k\backslash j,k\backslash j}^{jk}\hat{w}_{kj}.$$
Under Assumptions 3–5, S1 and S2,  $|\hat{\sigma}_{jk}^2 - \sigma_{jk}^2| = o_{\mathbb{P}}(1).$ 

 $\frac{1}{2} \frac{1}{2} \frac{1}$ 

*Proof of Theorem* 3. By Lemma S12,  $\hat{\sigma}_{jk}^2 = \sigma_{jk}^2 + o_{\mathbb{P}}(1)$ . The results follow from an application of Slutsky's theorem to (S29) in Lemma S10.

# S6. Proof of the lemmas in $\S$ S4

S6.1. Proof of Lemma S1

Proof. Recall from (S3) that

$$h_{ijk}^{rr'}(\beta_{j,O\setminus j}^{*}) = -\frac{R_{ij}^{rr'}(\beta_{j,O\setminus j}^{*})(x_{ij}^{r} - x_{ij}^{r'})(x_{ik}^{r} - x_{ik}^{rr'})}{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j}^{*})}.$$

By Proposition S1, for any t > 0 we have  $pr(|X_{ij}^r - X_{ij}^{r'}| \ge t) \le 2c_1 \exp(-t/2)$ . Note that  $R_{ij}^{rr'}(\beta_{j,O\setminus j}^*) > 0$  and therefore  $R_{ij}^{rr'}(\beta_{j,O\setminus j}^*)/\{1 + R_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\} < 1$ . Hence, by the union bound, we have that for any t > 0,

$$\operatorname{pr}\left\{h_{ijk}^{rr'}(\beta_{j,O\setminus j}^{*}) \ge t\right\} \leqslant \operatorname{pr}\left\{|(X_{ij}^{r} - X_{ij}^{r'})(X_{ik}^{r} - X_{ik}^{r'})| \ge t\right\}$$
  
$$\leqslant \operatorname{pr}\left(|X_{ij}^{r} - X_{ij}^{r'}| \ge t^{1/2}\right) + \operatorname{pr}\left(|X_{ik}^{r} - X_{ik}^{r'}| \ge t^{1/2}\right)$$
  
$$\leqslant 4c_{1} \exp(-t^{1/2}/2).$$

Then, by another application of the union bound,

$$\operatorname{pr}\left\{\frac{2}{R(R-1)}\sum_{1\leqslant r< r'\leqslant R}h_{ijk}^{rr'}(\beta_{j,O\setminus j}^*) \geqslant t\right\} = \operatorname{pr}\left\{\sum_{1\leqslant r< r'\leqslant R}h_{ijk}^{rr'}(\beta_{j,O\setminus j}^*) \geqslant tR(R-1)/2\right\}$$
$$\leqslant \sum_{1\leqslant r< r'\leqslant R}\operatorname{pr}\left\{h_{ijk}^{rr'}(\beta_{j,O\setminus j}^*) \geqslant t\right\}$$
$$\leqslant 4R^2c_1\exp(-t^{1/2}/2).$$

#### S6.2. Proof of Lemma S2

From (S1) and Lemma S1, we see that the gradient of the loss function is the average of <sup>285</sup> independent random variables with general exponential tail. To prove Lemma S2, we use a generalization of Bernstein's inequality for independent random variables with general exponential tail given in Lemma S14.

Recall from (S10) that  $U_{ik}^j(\beta_{j,O\setminus j}^*) = 2\sum_{r < r'} h_{ijk}^{rr'}(\beta_{j,O\setminus j}^*)/\{R(R-1)\}$ . By an application of Proposition 2,  $E\{U_{ik}^j(\beta_{j,O\setminus j}^*)\} = 0$ . In addition, by Lemma S1,

$$\operatorname{pr}\{U_{ik}^{j}(\beta_{j,O\setminus j}^{*}) \ge t\} \le 4c_1 R^2 \exp(-0.5t^{1/2}).$$

Since  $U_{1k}^j(\beta_{j,O\setminus j}^*), \ldots, U_{nk}^j(\beta_{j,O\setminus j}^*)$  are independent and identically distributed random variables with mean zero, by an application of Lemma S14 with  $L_1 = 4R^2c_1$ ,  $L_2 = 1/2$  and q = 1/2, we have

$$\Pr\left\{ \left| \nabla_k \ell_j(\beta_{j,O\setminus j}^*) \right| \ge t \right\} = \Pr\left\{ \left| \frac{1}{n} \sum_{i=1}^n U_{ik}^j(\beta_{j,O\setminus j}^*) \right| \ge t \right\}$$
  
$$\le 4 \exp\left( -\frac{1}{8} n^{1/5} t^{2/5} \right) + 16n R^2 c_1 \exp\left\{ -n^{1/5} t^{2/5} / (2^{3/2}) \right\}.$$

Therefore, by the union bound, we obtain

$$\Pr\{\|\nabla \ell_j(\beta_{j,O\setminus j}^*)\|_{\infty} \ge t\} \le 4p \exp\left(-\frac{1}{8}n^{1/5}t^{2/5}\right) + 16npR^2c_1 \exp\{-n^{1/5}t^{2/5}/(2^{3/2})\}.$$

Note that Lemma S14 holds only if  $t \ge (8E[\{U_{ik}^j(\beta_{j,O\setminus j}^*)\}^2]/n)^{1/2}$  for  $k \in O \setminus j$ . By Proposition S1,  $E[\{U_{ik}^j(\beta_{j,O\setminus j}^*)\}^2]$  is bounded since  $E(X_j^4)$  is bounded. We take  $t = K(\log^5 p/n)^{1/2}$  for sufficiently large K > 0 such that the inequality  $t \ge (8E[\{U_{ik}^j(\beta_{j,O\setminus j}^*)\}^2]/n)^{1/2}$  holds for all  $k \in O \setminus j$ . Then, for sufficiently large K, we have

$$\operatorname{pr}\left\{\left\|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\right\|_{\infty} \geqslant K\left(\frac{\log^{5} p}{n}\right)^{1/2}\right\} \leqslant p^{-1}.$$

We conclude that  $\|\nabla \ell_j(\beta_{j,O\setminus j}^*)\|_{\infty} \leq K(\log^5 p/n)^{1/2}$  with probability at least  $1-p^{-1}$ .

#### S6.3. Proof of Lemma S4

Recall that  $S_j$  is the support set of  $\beta_{j,O\setminus j}^*$  and that

$$D\left(\hat{\beta}_{j,O\setminus j},\beta_{j,O\setminus j}\right) = \left(\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j}\right)^{\mathrm{T}} \left\{ \nabla \ell_j(\hat{\beta}_{j,O\setminus j}) - \nabla \ell_j(\beta_{j,O\setminus j}) \right\}$$

is the symmetrized Bregman divergence of  $\ell_j(\beta_{j,O\setminus j})$ . Also, recall that  $\hat{\Delta}_j = \hat{\beta}_{j,O\setminus j} - \beta^*_{j,O\setminus j}$ . Observe that

$$D\left(\hat{\beta}_{j,O\setminus j},\beta_{j,O\setminus j}^{*}\right)$$
  
=  $\hat{\Delta}_{j}^{\mathrm{T}}\left\{\nabla\ell_{j}(\beta_{j,O\setminus j}^{*}+\hat{\Delta}_{j})-\nabla\ell_{j}(\beta_{j,O\setminus j}^{*})\right\}$   
=  $\sum_{k\in\mathcal{S}_{j}^{\mathrm{c}}}\hat{\beta}_{jk}\nabla_{k}\ell_{j}(\beta_{j,O\setminus j}^{*}+\hat{\Delta}_{j})+\sum_{k\in\mathcal{S}_{j}}\hat{\Delta}_{jk}\nabla_{k}\ell_{j}(\beta_{j,O\setminus j}^{*}+\hat{\Delta}_{j})-\hat{\Delta}_{j}^{\mathrm{T}}\nabla\ell_{j}(\beta_{j,O\setminus j}^{*}).$ 

290

By the Karush–Kuhn–Tucker conditions of (7),  $\hat{\beta}_{j,O\setminus j}$  is a solution to (7) if and only if

$$\begin{cases} \nabla_k \ell_j(\hat{\beta}_{j,O\setminus j}) = -\lambda \operatorname{sign}(\hat{\beta}_{jk}) & \text{if } \hat{\beta}_{jk} \neq 0, \\ \left| \nabla_k \ell_j(\hat{\beta}_{j,O\setminus j}) \right| \leqslant \lambda & \text{if } \hat{\beta}_{jk} = 0. \end{cases}$$

305

Note that 
$$k \in S_j^c$$
 does not imply  $\hat{\beta}_{jk} = 0$ , since  $S_j$  is the support set of  $\beta_{j,O\setminus j}^*$ .  
Nonetheless, by the Karush–Kuhn–Tucker conditions,  $\sum_{k \in S_j^c} \hat{\beta}_{jk} \nabla_k \ell_j (\beta_{j,O\setminus j}^* + \hat{\Delta}_j) = -\lambda \sum_{k \in S_j^c} \hat{\beta}_{jk} \operatorname{sign}(\hat{\beta}_{jk}) \text{ and } \sum_{k \in S_j} \hat{\Delta}_{jk} \nabla_k \ell_j (\beta_{j,O\setminus j}^* + \hat{\Delta}_j) \leq \lambda \sum_{k \in S_j} |\hat{\Delta}_{jk}|.$  Therefore  
 $D\left(\hat{\beta}_{j,O\setminus j}, \beta_{j,O\setminus j}^*\right) \leq -\lambda \sum_{k \in S_j^c} \hat{\beta}_{jk} \operatorname{sign}(\hat{\beta}_{jk}) + \lambda \sum_{k \in S_j} |\hat{\Delta}_{jk}| - \hat{\Delta}_j^T \nabla \ell_j (\beta_{j,O\setminus j}^*)$   
 $\leq -\lambda \|\hat{\Delta}_{j,S_j^c}\|_1 + \lambda \|\hat{\Delta}_{j,S_j}\|_1 + \|\hat{\Delta}_j\|_1 \|\nabla \ell_j (\beta_{j,O\setminus j}^*)\|_{\infty}$   
 $\leq -\lambda \|\hat{\Delta}_{j,S_j^c}\|_1 + \lambda \|\hat{\Delta}_{j,S_j}\|_1 + (\|\hat{\Delta}_{j,S_j}\|_1 + \|\hat{\Delta}_{j,S_j^c}\|_1) \|\nabla \ell_j (\beta_{j,O\setminus j}^*)\|_{\infty}$   
 $= \{\lambda + \|\nabla \ell_j (\beta_{j,O\setminus j}^*)\|_{\infty}\} \|\hat{\Delta}_{j,S_j}\|_1 - \{\lambda - \|\nabla \ell_j (\beta_{j,O\setminus j}^*)\|_{\infty}\} \|\hat{\Delta}_{j,S_j^c}\|_1.$ 
(S30)

The inequality is obtained by rearranging the terms in the last expression. To show that  $\|\hat{\Delta}_{j,\mathcal{S}_{j}^{c}}\|_{1} \leq \xi \|\hat{\Delta}_{j,\mathcal{S}_{j}}\|_{1}$ , we use the fact that  $D(\hat{\beta}_{j,O\setminus j},\beta_{j,O\setminus j}^{*}) \geq 0$  since the loss function  $\ell_{j}(\cdot)$  is a convex function. The inequality is obtained by substituting  $\|\nabla \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty} \leq \lambda(\xi-1)/(\xi+1)$  into (S30) and rearranging the terms.

# S6.4. Proof of Lemma S5

Recall the definitions of  $T_{ij}^{rr'}(\beta_{j,O\setminus j}^*)$ ,  $T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*)$  and  $\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)$  from (S4), (S5) and (S7), respectively. We first show that  $||E\{T_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\}||_{\infty}$  is bounded. By Proposition S1, for a sufficiently large constant C,

$$\left\| E\{T_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\} \right\|_{\infty} \leqslant C \max_{i,j,r} E|X_{ij}^r|^4 \leqslant 24Cc_1.$$

Hence, by the union bound, for any  $t \ge 2 \|E\{T_{ij}^{rr'}(\beta_{j,O\setminus j}^*)\}\|_{\infty}$ ,

$$\begin{aligned} \Pr\left[\left|T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^{*}) - E\{T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^{*})\}\right| \geqslant t\right] \\ &\leqslant \Pr\left[\left|T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^{*})\right| \geqslant t - \left|E\{T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^{*})\}\right|\right] \\ &\leqslant \Pr\left\{\left|T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^{*})\right| \geqslant t/2\right\} \\ &\leqslant \Pr\left\{\left|(X_{ij}^{r} - X_{ij}^{r'})^{2}(X_{ik}^{r} - X_{ik}^{r'})(X_{il}^{r} - X_{il}^{r'})\right| \geqslant t/2\right\} \\ &\leqslant \Pr\left\{\left|X_{ij}^{r} - X_{ij}^{r'}\right| \geqslant (t/2)^{1/4}\right\} + \Pr\left\{\left|X_{ik}^{r} - X_{ik}^{r'}\right| \geqslant (t/2)^{1/4}\right\} \\ &+ \Pr\left\{\left|X_{il}^{r} - X_{il}^{r'}\right| \geqslant (t/2)^{1/4}\right\} \\ &\leqslant 6c_{1}\exp\left(-t^{1/4}2^{-5/4}\right),\end{aligned}$$

where the last inequality follows from Proposition S1.

Choosing  $C_H = \max[6c_1, \exp\{(1.5c_1C)^{1/4}\}]$ , we have that for any t > 0,

$$\Pr\left[\left|T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*) - E\{T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*)\}\right| \ge t\right] \le C_H \exp\left(-t^{1/4}2^{-5/4}\right).$$

Let

$$V_{ikl} = \frac{2}{R(R-1)} \sum_{1 \le r < r' \le R} \left[ T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*) - E\{T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*)\} \right].$$

Then, by the union bound,

$$\operatorname{pr}\left(V_{ikl} \ge t\right) \leqslant \sum_{1 \leqslant r < r' \leqslant R} \operatorname{pr}\left[\left|T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*) - E\{T_{ijkl}^{rr'}(\beta_{j,O\setminus j}^*)\}\right| \ge t\right]$$
$$\leqslant C_H R^2 \exp\left(-t^{1/4} 2^{-5/4}\right).$$

By the definition of  $V_{ikl}$ , we have  $E(V_{ikl}) = 0$ . Since  $V_{1kl}, \ldots, V_{nkl}$  are independent and identically distributed random variables with mean zero, by an application of Lemma S14 with  $L_1 = C_H R^2$ ,  $L_2 = 2^{-5/4}$  and q = 1/4, we have that for any  $t \ge \{E(V_{ijk}^2)/n\}^{1/2}$ ,

$$\operatorname{pr}\left\{ \left| \nabla_{kl}^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) - H_{kl}^{j} \right| \ge t \right\} = \operatorname{pr}\left( \left| \frac{1}{n} \sum_{i=1}^{n} V_{ikl} \right| \ge t \right)$$
  
$$\le 4 \exp\left( -\frac{1}{8} n^{1/9} t^{2/9} \right) + 4n C_{H} R^{2} \exp\left\{ -n^{1/9} t^{2/9} / (2^{3/2}) \right\}.$$

Therefore, by the union bound, we obtain

$$\operatorname{pr}\left\{ \left\| \nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*}) - H^{j} \right\|_{\infty} \geq t \right\}$$

$$\leq 4p^{2} \exp\left(-\frac{1}{8}n^{1/9}t^{2/9}\right) + 4np^{2}C_{H}R^{2} \exp\left\{-n^{1/9}t^{2/9}/(2^{3/2})\right\}.$$

$$(S31)$$

Note that (S31) holds only if  $t \ge \{E(V_{ijk}^2)/n\}^{1/2}$ . It can be verified that  $E(V_{ijk}^2)$  is bounded since  $E(X_j^8)$  is bounded by Proposition S1. Therefore, taking  $t = K_1(\log^9 p/n)^{1/2}$  for sufficiently large  $K_1$ , we obtain

$$\operatorname{pr}\left(\left\|\nabla^{2}\ell_{j}(\beta_{j,O\setminus j}^{*})-H^{j}\right\|_{\infty} \geq t\right) \leq p^{-1}.$$

We conclude that  $\|\nabla^2 \ell_j(\beta^*_{j,O\setminus j}) - H^j\|_{\infty} \leq K_1 (\log^9 p/n)^{1/2}$  with probability at least  $1 - p^{-1}$ .

#### S7. Proof of the Lemmas in $\S$ S5

# S7.1. *Proof of Lemma* S6

Recall from §2.4 that the parameters associated with the *j*th and *k*th nodes are  $\beta_{j\vee k} = (\beta_{jk}, \beta_{j,O\setminus\{j,k\}}^{T}, \beta_{k,O\setminus\{j,k\}}^{T})^{T} \in \mathbb{R}^{2p-3}$ . Let  $L_{jk}(\beta_{j\vee k}) = \ell_{j}(\beta_{j,O\setminus j}) + \ell_{k}(\beta_{k,O\setminus k})$ . The gradient of  $L_{jk}(\beta_{j\vee k})$  evaluated at  $\beta_{j\vee k}$  is

$$\nabla_{jk}L_{jk}(\beta_{j\vee k}) = \frac{\partial L_{jk}(\beta_{j\vee k})}{\partial\beta_{jk}} = \nabla_k \ell_j(\beta_{j,O\setminus j}) + \nabla_j \ell_k(\beta_{k,O\setminus k}) \in \mathbb{R},$$
  
$$\nabla_{j,-k}L_{jk}(\beta_{j\vee k}) = \frac{\partial L_{jk}(\beta_{j\vee k})}{\partial\beta_{j,O\setminus\{j,k\}}} = \nabla_{-k}\ell_j(\beta_{j,O\setminus j}) \in \mathbb{R}^{p-2},$$
 (S32) 334

320

K. M. TAN, Y NING, D. M. WITTEN AND H. LIU

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

From (S32) and the definition of  $S_{jk}(\beta_{j\vee k})$  in (9), we see that  $S_{jk}(\beta_{j\vee k})$  is a linear transformation of  $\nabla L_{jk}(\beta_{j\vee k}) \in \mathbb{R}^{2p-3}$ . Let  $b = \{1, (-w_{jk}^*)^{\mathrm{T}}, (-w_{kj}^*)^{\mathrm{T}}\}^{\mathrm{T}} \in \mathbb{R}^{2p-3}$ . Then it can be verified that

$$n^{1/2}S_{jk}(\beta_{j\vee k}^{*}) = n^{1/2}b^{\mathrm{T}}\nabla L_{jk}(\beta_{j\vee k}^{*}) = \frac{1}{n^{1/2}}\sum_{i=1}^{n}b^{\mathrm{T}}g_{i}^{jk}(\beta_{j\vee k}^{*}),$$

where  $g_i^{jk}(\beta_{j\vee k}^*)$  is as defined in (14). By Proposition 2,  $g_1^{jk}(\beta_{j\vee k}^*), \ldots, g_n^{jk}(\beta_{j\vee k}^*)$  are independent and identically distributed random variables with mean zero and  $\operatorname{var}\{g_i^{jk}(\beta_{i\vee k}^*)\} = \Sigma^{jk}$ , where  $\Sigma^{jk}$  is as defined in (15). By an application of the central limit theorem, we have that  $n^{1/2}S_{jk}(\beta^*_{j\vee k})$  converges in distribution to  $N(0,\sigma^2_{jk})$ , where

$$\sigma_{jk}^2 = \Sigma_{jk,jk}^{jk} - 2\Sigma_{jk,j\backslash k}^{jk} w_{jk}^* - 2\Sigma_{jk,k\backslash j}^{jk} w_{kj}^* + (w_{jk}^*)^{\mathrm{T}} \Sigma_{j\backslash k,j\backslash k}^{jk} w_{jk}^* + (w_{kj}^*)^{\mathrm{T}} \Sigma_{k\backslash j,k\backslash j}^{jk} w_{kj}^*.$$

S7.2. Proof of Lemma S7

Recall from §2.4 that  $H^j = E\{\nabla^2 \ell_j(\beta_{j,O\setminus j}^*)\}$ . By the triangle inequality, 345

$$\begin{split} \left\| \nabla^2 \ell_j(\hat{\beta}_{j,O\setminus j}) - H^j \right\|_{\infty} &\leqslant \left\| \nabla^2 \ell_j(\hat{\beta}_{j,O\setminus j}) - \nabla^2 \ell_j(\beta^*_{j,O\setminus j}) \right\|_{\infty} + \left\| \nabla^2 \ell_j(\beta^*_{j,O\setminus j}) - H^j \right\|_{\infty} \\ &= I_1 + I_2. \end{split}$$

By Lemma S5, we have  $I_2 = \mathcal{O}_{\mathbb{P}}\{(\log^9 p/n)^{1/2}\} = o_{\mathbb{P}}(1)$ . By Lemma S3,

$$\nabla^2 \ell_j(\hat{\beta}_{j,O\setminus j}) - \nabla^2 \ell_j(\beta^*_{j,O\setminus j}) \leqslant \{\exp(b) - 1\} \nabla^2 \ell_j(\beta^*_{j,O\setminus j})$$

where

$$b = \max_{1 \leq i \leq n} \max_{1 \leq r < r' \leq R} \left| (x_{ij}^r - x_{ij}^{r'}) (\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j}^*)^{\mathrm{T}} (x_{i,O\setminus j}^r - x_{i,O\setminus j}^{r'}) \right|$$
  
$$\leq M \|\hat{\beta}_{j,O\setminus j} - \beta_{j,O\setminus j}^*\|_{1},$$

with M as defined in Theorem 1. By Assumption 4,  $||H^j||_{\infty} = \mathcal{O}(1)$ . Therefore, we obtain 350

$$I_{1} \leq |\exp(b) - 1| \|\nabla^{2} \ell_{j}(\beta_{j,O\setminus j}^{*})\|_{\infty}$$
  
$$\leq |\exp(b) - 1| (I_{2} + \|H^{j}\|_{\infty})$$
  
$$= \mathcal{O}_{\mathbb{P}}(|b|) \{o_{\mathbb{P}}(1) + \mathcal{O}(1)\}$$
  
$$= \mathcal{O}_{\mathbb{P}}(Ms_{\max}\lambda/\kappa_{\min}^{2}),$$

where the first inequality follows from Hölder's inequality, the second inequality follows from an application of the triangle inequality, and the last equality is obtained from an application of Theorem 1.

S7.3. Proof of Lemma S8

Recall from § 2.4 that  $(w_{jk}^*)^{\mathrm{T}} = (H_{k,-k}^j)^{\mathrm{T}} (H_{-k,-k}^j)^{-1}$ . By the triangle inequality and the def-355 inition of  $w_{ik}^*$ , we obtain

$$\left\|\nabla_{k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j})-(w_{jk}^{*})^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j})\right\|_{\infty}$$

Latent variable graphical models

$$\leq \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,O\setminus j}) - H_{k,-k}^{j} \right\|_{\infty} + \left\| (w_{jk}^{*})^{\mathsf{T}} \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,O\setminus j}) - H_{k,-k}^{j} \right\|_{\infty}$$

$$= \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,O\setminus j}) - H_{k,-k}^{j} \right\|_{\infty} + \left\| (w_{jk}^{*})^{\mathsf{T}} \left\{ \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,O\setminus j}) - H_{-k,-k}^{j} \right\} \right\|_{\infty}$$

$$= I_{1} + I_{2}.$$

By Lemma S7, we have  $I_1 = \mathcal{O}_{\mathbb{P}}\{Ms_{\max}\lambda/\kappa_{\min}^2 + (\log^9 p/n)^{1/2}\}$ . Similarly, by Hölder's inequality and Lemma S7, we have  $I_2 = \mathcal{O}_{\mathbb{P}}[w_0\{Ms_{\max}\lambda/\kappa_{\min}^2 + (\log^9 p/n)^{1/2}\}]$ . Therefore,

$$\begin{aligned} \left\| \nabla_{k,-k}^2 \ell_j(\hat{\beta}_{j,O\setminus j}) - (w_{jk}^*)^{\mathrm{T}} \nabla_{-k,-k}^2 \ell_j(\hat{\beta}_{j,O\setminus j}) \right\|_{\infty} \\ &= \mathcal{O}_{\mathbb{P}} \Big[ w_0 \Big\{ M s_{\max} \lambda / \kappa_{\min}^2 + (\log^9 p/n)^{1/2} \Big\} \Big]. \end{aligned}$$

Picking  $\lambda_w \ge Cw_0 \{Ms_{\max}\lambda/\kappa_{\min}^2 + (\log^9 p/n)^{1/2}\}$  for some sufficiently large C, we have

$$\left\|\nabla_{k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j})-(w_{jk}^{*})^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(\hat{\beta}_{j,O\setminus j})\right\|_{\infty} \leqslant \lambda_{w}.$$

# S7.4. Proof of Lemma S9

By Lemma S8, for  $\lambda_w \ge Cw_0 \{Ms_{\max}\lambda/\kappa_{\min}^2 + (\log^9 p/n)^{1/2}\}, w_{jk}^*$  is in the feasible region of the Dantzig selector problem (10); that is,

$$\left\|\nabla_{k,-k}^{2}\ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) - (w_{jk}^{*})^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}})\right\|_{\infty} \leq \lambda_{w}$$
(S33)

with high probability. For notational convenience, we let  $\hat{\Delta} = \hat{w}_{jk} - w_{jk}^*$ . By the triangle inequality,

$$\begin{split} \left\| \hat{\Delta}^{\mathrm{T}} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \right\|_{\infty} &\leq \left\| \nabla^{2}_{k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) - \hat{w}^{\mathrm{T}}_{jk} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \right\|_{\infty} \\ &+ \left\| \nabla^{2}_{k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) - (w^{*}_{jk})^{\mathrm{T}} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \right\|_{\infty} \\ &\leq 2\lambda_{w}, \end{split}$$
(S34)

where the last inequality follows from (S33) and the constraint in (10). By Hölder's inequality and (S34),

$$\hat{\Delta}^{\mathrm{T}} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \hat{\Delta} \leqslant \|\hat{\Delta}\|_{1} \|\hat{\Delta}^{\mathrm{T}} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}})\|_{\infty} \leqslant 2\lambda_{w} \|\hat{\Delta}\|_{1}.$$
 (S35)

Let  $\mathcal{B}_{jk}$  be the support set of  $w_{jk}^*$ , i.e.,  $\mathcal{B}_{jk} = \{l : (w_{jk}^*)_l \neq 0\}$ . Also, let  $s'_{jk} = |\mathcal{B}_{jk}|$  be the cardinality of  $\mathcal{B}_{jk}$ . By the definition of the Dantzig selector, (10),  $\|\hat{w}_{jk}\|_1 \leq \|w_{jk}^*\|_1$ , implying

$$\sum_{l \in \mathcal{B}_{jk}} |(w_{jk}^*)_l| \ge \sum_{l \in \mathcal{B}_{jk}} |(\hat{w}_{jk})_l| + \sum_{l \in \mathcal{B}_{jk}^c} |(\hat{w}_{jk})_l|.$$
 (S36)

By the triangle inequality,

$$\sum_{l \in \mathcal{B}_{jk}} |(\hat{w}_{jk})_l - (w_{jk}^*)_l| \ge \sum_{l \in \mathcal{B}_{jk}} |(w_{jk}^*)_l| - \sum_{l \in \mathcal{B}_{jk}} |(\hat{w}_{jk})_l|.$$
(S37)

Upon adding (S36) and (S37) and rearranging the terms, we obtain  $\|\hat{\Delta}_{\mathcal{B}_{jk}^c}\|_1 \leq \|\hat{\Delta}_{\mathcal{B}_{jk}}\|_1$ , which implies

$$\|\hat{\Delta}\|_1 \leqslant 2 \|\hat{\Delta}_{\mathcal{B}_{ik}}\|_1. \tag{S38}$$

360

370

Substituting (S38) into (S35) gives

$$\hat{\Delta}^{\mathrm{T}} \nabla^2_{-k,-k} \ell_j(0, \hat{\beta}_{j,O\setminus\{j,k\}}) \hat{\Delta} \leqslant 4\lambda_w \| \hat{\Delta}_{\mathcal{B}_{jk}} \|_1.$$
(S39)

We now derive a lower bound for  $\hat{\Delta}^{\mathrm{T}} \nabla^2_{-k,-k} \ell_j(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \hat{\Delta}$ . Let  $\mathcal{A} = (v \in \mathbb{R}^{p-2} : \|v_{\mathcal{B}_{jk}^c}\|_1 \leq \|v_{\mathcal{B}_{jk}}\|_1)$ . We first show that for sufficiently large n, the quantity

$$\inf_{v \in \mathcal{A}} \frac{s_{jk}' \{ v^{\mathrm{T}} \nabla_{-k,-k}^2 \ell_j(0,\hat{\beta}_{j,O\setminus\{j,k\}}) v \}}{\| v_{\mathcal{B}_{jk}} \|_1^2}$$

can be bounded below by a positive constant. By an application of Lemma S3, we obtain

$$\frac{s_{jk}^{\prime}\left\{v^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}})v\right\}}{\|v_{\mathcal{B}_{jk}}\|_{1}^{2}} \geqslant \frac{s_{jk}^{\prime}\left\{v^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})v\right\}}{\|v_{\mathcal{B}_{jk}}\|_{1}^{2}}\exp(-b), \quad (S40)$$

380 where

$$b = \max_{1 \le i \le n} \max_{1 \le r < r' \le R} \left| (x_{ij}^r - x_{ij}^{r'}) (0, \hat{\beta}_{j,O \setminus \{j,k\}} - \beta_{j,O \setminus \{j,k\}}^*)^{\mathrm{T}} (x_{i,O \setminus j}^r - x_{i,O \setminus j}^{r'}) \right|.$$

By Theorem 1 and Assumption S1,

$$b \leqslant M \|\hat{\beta}_{j,O\setminus\{j,k\}} - \beta_{j,O\setminus\{j,k\}}^*\|_1 \leqslant \frac{M \exp(\eta)(\xi+1)}{2\kappa_{\min}^2} s_{\max}\lambda \leqslant \log 2$$
(S41)

for sufficiently large n.

By (S40), (S41) and Assumption 4 that  $\Lambda_{\min}(H^j) \ge \Lambda_{\text{lower}}^H > 0$ , we have

$$\frac{s_{jk}'\{v^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}})v\}}{\|v_{\mathcal{B}_{jk}}\|_{1}^{2}}} \\
\geq \frac{1}{2} \frac{s_{jk}'\{v^{\mathrm{T}}\nabla_{-k,-k}^{2}\ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})v\}}{\|v_{\mathcal{B}_{jk}}\|_{1}^{2}}} \\
= \frac{s_{jk}'v^{\mathrm{T}}\{H_{-k,-k}^{j} - H_{-k,-k}^{j} + \nabla_{-k,-k}^{2}\ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})\}v}{2\|v_{\mathcal{B}_{jk}}\|_{1}^{2}} \\
\geq \frac{s_{jk}'\|v\|_{2}^{2}\Lambda_{\mathrm{lower}}^{H} - s_{jk}'\|H_{-k,-k}^{j} - \nabla_{-k,-k}^{2}\ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})\|_{\infty}\|v\|_{1}^{2}}{2\|v_{\mathcal{B}_{jk}}\|_{1}^{2}},$$
(S42)

where the last expression follows by an application of Hölder's inequality. Noting that  $||v_{\mathcal{B}_{jk}}||_1 \leq (s'_{jk})^{1/2} ||v_{\mathcal{B}_{jk}}||_2 \leq (s'_{jk})^{1/2} ||v||_2$  and  $||v||_1^2 \leq 4 ||v_{\mathcal{B}_{jk}}||_1^2$  for any  $v \in \mathcal{A}$ , we obtain

$$\frac{s'_{jk} \|v\|_{2}^{2} \Lambda_{\text{lower}}^{H} - s'_{jk} \|H_{-k,-k}^{j} - \nabla_{-k,-k}^{2} \ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})\|_{\infty} \|v\|_{1}^{2}}{2\|v_{\mathcal{B}_{jk}}\|_{1}^{2}} \\
\geqslant \frac{1}{2} \Lambda_{\text{lower}}^{H} - 2s'_{jk} \|H_{-k,-k}^{j} - \nabla_{-k,-k}^{2} \ell_{j}(0,\beta_{j,O\setminus\{j,k\}}^{*})\|_{\infty} \\
\geqslant \frac{1}{2} \Lambda_{\text{lower}}^{H} - 2K_{1}s'_{jk} (\log^{9} p/n)^{1/2} \\
= \frac{1}{2} \Lambda_{\text{lower}}^{H} + o_{\mathbb{P}}(1),$$
(S43)

where the last inequality is obtained by applying Lemma S7 and the last equality follows from Assumption S1.

Upon setting  $v = \hat{\Delta}$  and combining (S42) and (S43), we obtain

$$\hat{\Delta}^{\mathrm{T}} \nabla^{2}_{-k,-k} \ell_{j}(0,\hat{\beta}_{j,O\setminus\{j,k\}}) \hat{\Delta} \geqslant \frac{1}{2} \Lambda^{H}_{\mathrm{lower}}(s'_{jk})^{-1} \|\hat{\Delta}_{\mathcal{B}_{jk}}\|_{1}^{2}.$$
(S44)

Finally, combining (S39) and (S44) yields

$$\|\hat{\Delta}_{\mathcal{B}_{jk}}\|_1 \leqslant 8s'_{jk}\lambda_w / \Lambda^H_{\text{lower}},$$

which implies

$$\|\hat{w}_{jk} - w_{jk}^*\|_1 = \|\hat{\Delta}\|_1 \leqslant 2\|\hat{\Delta}_{\mathcal{B}_{jk}}\|_1 \leqslant 16s'_{jk}\lambda_w/\Lambda^H_{\text{lower}} = \mathcal{O}_{\mathbb{P}}(s'_{\max}\lambda_w),$$

where the first inequality is obtained from (S38).

#### S7.5. Proof of Lemma S10

The proof consists of two parts. We first show that under the null hypothesis  $H_0: \beta_{jk}^* = \beta_{kj}^* = 0$  in (8),  $n^{1/2}\hat{S}_{jk} = n^{1/2}S_{jk}(\beta_{j\vee k}^*) + o_{\mathbb{P}}(1)$ . Then, by an application of Lemma S6, we show that  $n^{1/2}\hat{S}_{jk}$  is asymptotically normal.

Recall the definitions of  $S_{jk}(\beta_{j\vee k}^*)$  and  $\hat{S}_{jk}$  from (9) and (11), respectively. With some abuse of notation, throughout this proof we write  $\hat{\beta}_{j,-k} = (0, \hat{\beta}_{j,O\setminus\{j,k\}})$  and  $\beta_{j,-k}^* = (0, \beta_{j,O\setminus\{j,k\}}^*)$ . Under the null hypothesis  $H_0: \beta_{jk}^* = \beta_{kj}^* = 0$ , we have that  $\hat{S}_{jk} - S_{jk}(\beta_{j\vee k}^*) = I_{1j} + I_{2j} + I_{1k} + I_{2k}$ , where  $I_{1j}$  and  $I_{2j}$  are defined as

$$I_{1j} = \nabla_k \ell_j(\hat{\beta}_{j,-k}) - \nabla_k \ell_j(\beta_{j,-k}^*) - \hat{w}_{jk}^{\mathrm{T}} \{ \nabla_{-k} \ell_j(\hat{\beta}_{j,-k}) - \nabla_{-k} \ell_j(\beta_{j,-k}^*) \}, \qquad 400$$

$$I_{2j} = (w_{jk}^* - \hat{w}_{jk})^{\mathrm{T}} \nabla_{-k} \ell_j(\beta_{j,-k}^*).$$

The terms  $I_{1k}$  and  $I_{2k}$  are defined similarly by interchanging the subscripts j and k in  $I_{1j}$  and  $I_{2j}$ . The goal is to show that each of the four terms is  $o_{\mathbb{P}}(n^{-1/2})$ .

Upper bound for  $I_{1j}$ . Let  $\hat{\Delta}_{j,-k} = \hat{\beta}_{j,-k} - \beta^*_{j,-k}$ . By an application of the mean value theorem, there exists  $\tilde{\beta}_{j,-k} \in \mathbb{R}^{p-1}$  on the line segment between  $\hat{\beta}_{j,-k}$  and  $\beta^*_{j,-k}$  such that

$$I_{1j} = \left\{ \nabla_{k,-k}^2 \ell_j(\tilde{\beta}_{j,-k}) - \hat{w}_{jk}^{\mathrm{T}} \nabla_{-k,-k}^2 \ell_j(\tilde{\beta}_{j,-k}) \right\} \hat{\Delta}_{j,-k}$$

By the triangle inequality and Hölder's inequality, we have

$$\begin{split} |I_{1j}| &\leqslant \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \hat{w}_{jk}^{\mathrm{T}} \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &+ \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &+ \left\| \hat{w}_{jk}^{\mathrm{T}} \left\{ \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{-k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) \right\} \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &= I_{1j1} + I_{1j2} + I_{1j3}. \end{split}$$

We now obtain an upper bound for each of the three terms separately.

By Theorem 1, the definition (10) of a Dantzig selector-type estimator, and Assumption S1, we have

$$I_{1j1} = \mathcal{O}_{\mathbb{P}}\left(\frac{s_{\max}\lambda\lambda_w}{\kappa_{\min}^2}\right) = o_{\mathbb{P}}(n^{-1/2}).$$

395

405

<sup>410</sup> By the triangle inequality,

$$\begin{split} I_{1j2} &= \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &\leq \left\| \nabla_{k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{k,-k}^{2} \ell_{j}(\beta_{j,-k}^{*}) \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &+ \left\| \nabla_{k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) - \nabla_{k,-k}^{2} \ell_{j}(\beta_{j,-k}^{*}) \right\|_{\infty} \|\hat{\Delta}_{j,-k}\|_{1} \\ &= \mathcal{O}_{\mathbb{P}}\left( \frac{Ms_{\max}\lambda}{\kappa_{\min}^{2}} \right) \mathcal{O}_{\mathbb{P}}\left( \frac{s_{\max}\lambda}{\kappa_{\min}^{2}} \right), \end{split}$$

where the last equality follows from the proof of Lemma S7 and Theorem 1. We can write the last expression as  $I_{ij2} = \mathcal{O}_{\mathbb{P}}(s_{\max}\lambda\lambda_w/\kappa_{\min}^2) = o_{\mathbb{P}}(n^{-1/2})$ , since  $\lambda_w = \Omega\{w_0(Ms_{\max}\lambda/\kappa_{\min}^2 + \lambda \log^2 p)\}$  by Assumption S1.

By the definition (10) of the Dantzig selector-type estimator,  $\|\hat{w}_{jk}\|_1 \leq \|w_{jk}^*\|_1$ . Therefore, by the triangle inequality and Hölder's inequality, we obtain

$$\begin{split} I_{1j3} &= \left\| \hat{w}_{jk}^{\mathrm{T}} \left\{ \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{-k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) \right\} \right\|_{\infty} \| \hat{\Delta}_{j,-k} \|_{1} \\ &\leq \| w_{jk}^{*} \|_{1} \left\| \nabla_{-k,-k}^{2} \ell_{j}(\hat{\beta}_{j,-k}) - \nabla_{-k,-k}^{2} \ell_{j}(\tilde{\beta}_{j,-k}) \right\|_{\infty} \| \hat{\Delta}_{j,-k} \|_{1} \\ &= \| w_{jk}^{*} \|_{1} \mathcal{O}_{\mathbb{P}} \left( \frac{M s_{\max} \lambda}{\kappa_{\min}^{2}} \right) \mathcal{O}_{\mathbb{P}} \left( \frac{s_{\max} \lambda}{\kappa_{\min}^{2}} \right) \\ &= \mathcal{O}_{\mathbb{P}} \left( \frac{s_{\max} \lambda \lambda_{w}}{\kappa_{\min}^{2}} \right) \\ &= o_{\mathbb{P}} (n^{-1/2}), \end{split}$$

420

where the equalities hold by Assumption S1.

Upper bound for  $I_{2j}$ . By Lemmas S2–S9 and an application of the Hölder's inequality, we obtain

$$|I_{2j}| \leq \|w_{jk}^* - \hat{w}_{jk}\|_1 \|\nabla_{-k}\ell_j(\beta_{j,-k}^*)\|_\infty = \mathcal{O}_{\mathbb{P}}(s_{\max}'\lambda_w)\mathcal{O}_{\mathbb{P}}(\lambda) = o_{\mathbb{P}}(n^{-1/2}),$$

where the last equality holds by Assumption S1.

<sup>425</sup> Combining the upper bounds for  $I_{1j}$  and  $I_{2j}$ , we obtain  $I_{1j} + I_{2j} = o_{\mathbb{P}}(n^{-1/2})$ . Similarly,  $I_{1k} + I_{2k} = o_{\mathbb{P}}(n^{-1/2})$ . Thus,  $n^{1/2}\{\hat{S}_{jk} - S_{jk}(\beta_{j\vee k}^*)\} = o_{\mathbb{P}}(1)$ . Asymptotic normality of  $n^{1/2}\hat{S}_{jk}$  is established by an application of Lemma S6.

# S7.6. Proof of Lemma S11

Recall from (14) and (15) the definitions of  $g_i^{jk}(\beta_{j\vee k})$  and  $\Sigma^{jk}$ , respectively. Also, recall from (16) that

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

For notational convenience, we write  $\hat{\beta}'_{j\vee k} = (0, \hat{\beta}_{j,O\setminus\{j,k\}}, \hat{\beta}_{k,O\setminus\{j,k\}})$ . By the triangle inequality, we obtain

$$\begin{aligned} \left\| \hat{\Sigma}^{jk}(\hat{\beta}'_{j\vee k}) - \Sigma^{jk} \right\|_{\infty} &\leq \left\| \hat{\Sigma}^{jk}(\beta^*_{j\vee k}) - \Sigma^{jk} \right\|_{\infty} + \left\| \hat{\Sigma}^{jk}(\hat{\beta}'_{j\vee k}) - \hat{\Sigma}^{jk}(\beta^*_{j\vee k}) \right\|_{\infty} \\ &= I_1 + I_2. \end{aligned}$$

We now obtain an upper bound for each term separately.

Upper bound for  $I_1$ . For notational simplicity, throughout the proof we suppress the superscript jk and write  $g_i^*$  to indicate  $g_i^{jk}(\beta_{j\vee k}^*)$ . We also write  $U_{ik}^j$  for  $U_{ik}^j(\beta_{j,O\setminus j}^*)$ , where  $U_{ik}^j(\beta_{i,O\setminus j}^*)$  is as defined in (S10). Note that

$$\hat{\Sigma}^{jk}(\beta_{j\vee k}^*) - \Sigma^{jk} = \frac{1}{n} \sum_{i=1}^n \left[ (g_i^*)^{\otimes 2} - E\{(g_i^*)^{\otimes 2}\} \right],$$

where  $(g_i^*)^{\otimes 2}$  takes the form

$$(g_{i}^{*})^{\otimes 2} = \begin{cases} (U_{ik}^{j} + U_{ij}^{k})^{2} & (U_{ik}^{j} + U_{ij}^{k})(U_{i,-k}^{j})^{\mathrm{T}} & (U_{ik}^{j} + U_{ij}^{k})(U_{i,-j}^{k})^{\mathrm{T}} \\ U_{i,-k}^{j}(U_{ik}^{j} + U_{ij}^{k}) & (U_{i,-k}^{j})^{\otimes 2} & U_{i,-k}^{j}(U_{i,-j}^{k})^{\mathrm{T}} \\ U_{i,-j}^{k}(U_{ik}^{j} + U_{ij}^{k}) & U_{i,-j}^{k}(U_{i,-k}^{j})^{\mathrm{T}} & (U_{i,-j}^{k})^{\otimes 2} \end{cases} \in \mathbb{R}^{(2p-3) \times (2p-3)}.$$

$$(S45)$$

We see from (S45) that each term in  $(g_i^*)^{\otimes 2} - E\{(g_i^*)^{\otimes 2}\}$  involves the quantity

$$U_{il}^j U_{im}^k - E(U_{il}^j U_{im}^k) \tag{S46}$$

for  $j \neq l$ ,  $k \neq m$ , and  $j, k, l, m \in O$ .

We now show that (S46) is bounded with high probability. First, note that  $||E(U_i^j U_i^k)||_{\infty} \leq C_u$  for some sufficiently large constant  $C_u$ , since  $E(X_j^4)$  is bounded by Proposition S1. Hence, by the union bound, for any  $t > 2||E(U_i^j U_i^k)||_{\infty}$  we have

$$\begin{split} \Pr\left\{ \left| U_{il}^{j} U_{im}^{k} - E(U_{il}^{j} U_{im}^{k}) \right| \geqslant t \right\} &\leq \Pr\left( \left| U_{il}^{j} U_{im}^{k} \right| \geqslant t/2 \right) \\ &\leq \Pr\{ |U_{il}^{j}| \geqslant (t/2)^{1/2} \} + \Pr\{ |U_{im}^{k}| \geqslant (t/2)^{1/2} \} \\ &\leq 8R^{2}c_{1} \exp(-t^{1/4}2^{-5/4}), \end{split}$$

where the last inequality is obtained by an application of Lemma S1. Note that the above inequality holds only if  $t > 2C_u$ . Choosing  $C_H = \max\{8R^2c_1, \exp(2^{-1}C_u^{1/4})\}$ , we have

$$\Pr\left\{\left|U_{il}^{j}U_{im}^{k} - E(U_{il}^{j}U_{im}^{k})\right| \ge t\right\} \le C_{H}\exp(-t^{1/4}2^{-5/4})$$

for all t > 0. Therefore, for any  $j \neq l$ ,  $k \neq m$  and  $j, k, l, m \in O$ , we have

$$\Pr\left\{ \left| \hat{\Sigma}^{jk}(\beta_{j\vee k}^{*}) - \Sigma^{jk} \right|_{jl,km} \ge t \right\}$$

$$\leq 4 \Pr\left[ \left| \frac{1}{n} \sum_{i=1}^{n} \left\{ U_{il}^{j} U_{im}^{k} - E(U_{il}^{j} U_{im}^{k}) \right\} \right| \ge \frac{t}{4} \right]$$

$$\leq 16 \exp\left( -\frac{1}{8} 4^{-2/9} n^{1/9} t^{2/9} \right) + 16n C_{H} \exp\left( -\frac{1}{8} 4^{-2/9} n^{1/9} t^{2/9} \right),$$

where the first inequality is obtained by an application of the union bound and the last inequality is obtained by an application of Lemma S14 with  $L_1 = C_H$ ,  $L_2 = 2^{-5/4}$  and q = 1/4. Hence, by the union bound, we obtain

$$\operatorname{pr}\left(\left\|\hat{\Sigma}^{jk}(\beta_{j\vee k}^{*}) - \Sigma^{jk}\right\|_{\infty} \ge t\right)$$

K. M. TAN, Y NING, D. M. WITTEN AND H. LIU

450

24

$$\leq 64p^{2} \exp\left(-\frac{1}{8}4^{-2/9}n^{1/9}t^{2/9}\right) + 64np^{2}C_{H} \exp\left(-\frac{1}{8}4^{-2/9}n^{1/9}t^{2/9}\right).$$

The above inequality holds only if  $t \ge [E\{U_{il}^j U_{im}^k - E(U_{il}^j U_{im}^k)\}^2/n]^{1/2}$ . By Proposition S1, the numerator is bounded since  $E(X_j^8)$  is bounded. Taking  $t = K_2(\log^9 p/n)^{1/2}$  for sufficiently large  $K_2 > 0$ , we have  $pr\{\|\hat{\Sigma}^{jk}(\beta_{j\vee k}^*) - \Sigma^{jk}\|_{\infty} \ge t\} \le p^{-1}$ . We conclude that

$$I_1 = \left\| \hat{\Sigma}^{jk}(\beta_{j\vee k}^*) - \Sigma^{jk} \right\|_{\infty} = \mathcal{O}_{\mathbb{P}}\left\{ (\log^9 p/n)^{1/2} \right\}.$$

Upper bound for  $I_2$ . For notational convenience, we let  $G^i(\beta_{j\vee k}) = \{g_i^{jk}(\beta_{j\vee k})\}^{\otimes 2} \in \mathbb{R}^{(2p-3)\times(2p-3)}$ . From (S45), we see that for any  $(a,b), (c,d) \in \{(p,q): p,q \in O, p \neq q\}$ ,  $G^i_{ab,cd}(\beta_{j\vee k})$  is a linear combination of  $U^j_{ik}(\beta_{j,O\setminus k})U^k_{ij}(\beta_{k,O\setminus j})$  and  $\{U^j_{ik}(\beta_{j,O\setminus k})\}^2$ . It can be shown that  $\partial U^j_{ik}(\beta_{j,O\setminus k})U^k_{ij}(\beta_{k,O\setminus j})/\partial \beta_{jk}$  involves the term  $(X^r_{ij} - X^{r'}_{ij})^3(X^r_{ik} - X^{r'}_{ik})^3$ . Since  $\max_{i,j,r} X^r_{ij} = \mathcal{O}_{\mathbb{P}}(\log p)$  by Proposition S1, there exists a constant  $C_v$  such that

$$\|\nabla G^i_{ab,cd}(\beta_{j\vee k})\|_{\infty} \leqslant C_v \log^6 p$$

for all  $j, k \in O$  with  $j \neq k$  and for any  $(a, b), (c, d) \in \{(p, q) : p, q \in O, p \neq q\}$ , with high probability. By the mean value theorem, there exists a  $\tilde{\beta}_{j\vee k}$  on the line segment between  $\hat{\beta}_{j\vee k}$  and  $\beta_{j\vee k}^*$  such that

$$\begin{aligned} G^{i}_{ab,cd}(\hat{\beta}_{j\vee k}) - G^{i}_{ab,cd}(\beta^{*}_{j\vee k}) &= \nabla G^{i}_{ab,cd}(\tilde{\beta}_{j\vee k})(\hat{\beta}_{j\vee k} - \beta^{*}_{j\vee k}) \\ &\leqslant \left\| \nabla G^{i}_{ab,cd}(\tilde{\beta}_{j\vee k}) \right\|_{\infty} \left\| \hat{\beta}_{j\vee k} - \beta^{*}_{j\vee k} \right\|_{1} \\ &= \mathcal{O}_{\mathbb{P}}(s_{\max}\lambda \log^{6} p/\kappa_{\min}^{2}), \end{aligned}$$

<sup>465</sup> where the last equality holds by an application of Theorem 1. Thus,

$$||I_2||_{\infty} = \mathcal{O}_{\mathbb{P}}(s_{\max}\lambda\log^6 p/\kappa_{\min}^2).$$

The result is obtained by combining the upper bounds for  $I_1$  and  $I_2$ .

S7.7. Proof of Lemma S12

Recall from (15) and (16) the definitions of  $\Sigma^{jk}$  and  $\hat{\Sigma}^{jk}$ , respectively. Let

$$\sigma_{jk}^2 = \Sigma_{jk,jk}^{jk} - 2\Sigma_{jk,j\backslash k}^{jk} w_{jk}^* - 2\Sigma_{jk,k\backslash j}^{jk} w_{kj}^* + (w_{jk}^*)^{\mathrm{T}} \Sigma_{j\backslash k,j\backslash k}^{jk} w_{jk}^* + (w_{kj}^*)^{\mathrm{T}} \Sigma_{k\backslash j,k\backslash j}^{jk} w_{kj}^*$$

and

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

470

To show that  $|\hat{\sigma}_{jk}^2 - \sigma_{jk}^2| = o_{\mathbb{P}}(1)$ , we use the results from Lemmas S9 and S11 and the fact that  $\|\Sigma^{jk}\|_{\infty} = \mathcal{O}(1)$  since  $E(X_j^4)$  is bounded by Proposition S1. For notational simplicity, we suppress the superscripts in  $\hat{\Sigma}^{jk}$  and  $\Sigma^{jk}$  in the following proof.

By the triangle inequality,

$$\begin{aligned} \left| \hat{\sigma}_{jk}^{2} - \sigma_{jk}^{2} \right| &\leq \left| \hat{\Sigma}_{jk,jk} - \Sigma_{jk,jk} \right| + 2 \left| \hat{\Sigma}_{jk,j\backslash k} \hat{w}_{jk} - \Sigma_{jk,j\backslash k} w_{jk}^{*} \right| + 2 \left| \hat{\Sigma}_{jk,k\backslash j} \hat{w}_{kj} - \Sigma_{jk,k\backslash j} w_{kj}^{*} \right| \\ &+ \left| \hat{w}_{jk}^{\mathrm{T}} \hat{\Sigma}_{j\backslash k,j\backslash k} \hat{w}_{jk} - (w_{jk}^{*})^{\mathrm{T}} \Sigma_{j\backslash k,j\backslash k} w_{jk}^{*} \right| \\ &+ \left| \hat{w}_{kj}^{\mathrm{T}} \hat{\Sigma}_{k\backslash j,k\backslash j} \hat{w}_{kj} - (w_{kj}^{*})^{\mathrm{T}} \Sigma_{k\backslash j,k\backslash j} w_{kj}^{*} \right| \end{aligned}$$

$$= I_1 + I_{2j} + I_{2k} + I_{3j} + I_{3k}$$

Upper bound for  $I_1$ . By Lemma S11, we obtain

$$I_1 = \mathcal{O}_{\mathbb{P}}\left\{\frac{s_{\max}}{\kappa_{\min}^2}\,\lambda\log^6 p + \left(\frac{\log^9 p}{n}\right)^{1/2}\right\} = \mathcal{O}_{\mathbb{P}}\left(\frac{s_{\max}}{\kappa_{\min}^2}\,\lambda\log^6 p\right).$$

Upper bound for  $I_{2j}$ . By the triangle inequality,

$$\begin{split} I_{2j} &= \left| \hat{\Sigma}_{jk,j \setminus k} \hat{w}_{jk} - \Sigma_{jk,j \setminus k} w_{jk}^* \right| \\ &\leq \left| (\hat{\Sigma}_{jk,j \setminus k} - \Sigma_{jk,j \setminus k}) (\hat{w}_{jk} - w_{jk}^*) \right| + \left| \Sigma_{jk,j \setminus k} (\hat{w}_{jk} - w_{jk}^*) \right| + \left| (\hat{\Sigma}_{jk,j \setminus k} - \Sigma_{jk,j \setminus k}) w_{jk}^* \right| \\ &= I_{2j1} + I_{2j2} + I_{2j3}. \end{split}$$

By Hölder's inequality and Lemmas S9 and S11, we have

$$I_{2j1} \leq \|\hat{w}_{jk} - w_{jk}^*\|_1 \|\Sigma_{jk,j\setminus k} - \Sigma_{jk,j\setminus k}\|_{\infty}$$
  
=  $\mathcal{O}_{\mathbb{P}}(s'_{\max}\lambda_w)\mathcal{O}_{\mathbb{P}}\{s_{\max}\lambda\log^6 p/\kappa_{\min}^2 + (\log^9 p/n)^{1/2}\}$   
=  $\mathcal{O}_{\mathbb{P}}(s'_{\max}\lambda_w s_{\max}\lambda\log^6 p/\kappa_{\min}^2).$ 

By Hölder's inequality, Lemma S9, and the fact that  $\|\Sigma\|_{\infty} = \mathcal{O}(1)$ , we obtain

$$I_{2j2} \leqslant \|\Sigma_{jk,j\setminus k}\|_{\infty} \|\hat{w}_{jk} - w_{jk}^*\|_1 = \mathcal{O}_{\mathbb{P}}(s_{\max}'\lambda_w).$$

Recall that  $w_0 = \max_{j,k \in O} ||w_{jk}^*||_1$ . By Hölder's inequality and Lemma S11,

$$I_{2j3} \leqslant \|\hat{\Sigma}_{jk,j\setminus k} - \Sigma_{jk,j\setminus k}\|_{\infty} \|w_{jk}^*\|_1 = \mathcal{O}_{\mathbb{P}}(w_0 s_{\max} \lambda \log^6 p/\kappa_{\min}^2).$$

Combining the upper bounds, we obtain

$$I_{2j} = \mathcal{O}_{\mathbb{P}} \{ (s'_{\max}\lambda_w + w_0) s_{\max}\lambda \log^6 p / \kappa_{\min}^2 + s'_{\max}\lambda_w \}$$
  
=  $\mathcal{O}_{\mathbb{P}} [\{o(1) + w_0\} s_{\max}\lambda \log^6 p / \kappa_{\min}^2 + o(1)]$   
=  $\mathcal{O}_{\mathbb{P}} (w_0 s_{\max}\lambda \log^6 p / \kappa_{\min}^2),$ 

where we have used the fact that  $s'_{\max}\lambda_w = o(1)$  by Assumption S1.

Upper bound for  $I_3$ . By the triangle inequality,

$$\begin{split} I_{3j} &\leqslant \left| \hat{w}_{jk}^{\mathrm{T}} (\hat{\Sigma}_{j \setminus k, j \setminus k} - \Sigma_{j \setminus k, j \setminus k}) \hat{w}_{jk} \right| + \left| \hat{w}_{jk}^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} \hat{w}_{jk} - (w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} w_{jk}^{*} \right| \\ &= I_{3j1} + I_{3j2}. \end{split}$$

By the definition of  $\hat{w}_{jk}$  in (10), we have  $\|\hat{w}_{jk}\|_1 \leq \|w_{jk}^*\|_1$ . Therefore, by Hölder's inequality and Lemma S11,

$$I_{3j1} \leq \|w_{jk}^*\|_1^2 \|\Sigma_{j \setminus k, j \setminus k} - \Sigma_{j \setminus k, j \setminus k}\|_{\infty}$$
  
=  $\mathcal{O}_{\mathbb{P}} \Big[ w_0^2 \Big\{ s_{\max\lambda} \log^6 p / \kappa_{\min}^2 + (\log^9 p / n)^{1/2} \Big\} \Big]$   
=  $\mathcal{O}_{\mathbb{P}} \Big( w_0^2 s_{\max\lambda} \log^6 p / \kappa_{\min}^2 \Big).$ 

485

480

By the triangle inequality,

~

$$\begin{split} I_{3j2} &= \left| \hat{w}_{jk}^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} \hat{w}_{jk} - (w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} w_{jk}^{*} \right| \\ &= \left| (\hat{w}_{jk} - w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} (\hat{w}_{jk} - w_{jk}^{*}) + 2(w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} (\hat{w}_{jk} - w_{jk}^{*}) \right| \\ &\leq \left| (\hat{w}_{jk} - w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} (\hat{w}_{jk} - w_{jk}^{*}) \right| + 2 \left| (w_{jk}^{*})^{\mathrm{T}} \Sigma_{j \setminus k, j \setminus k} (\hat{w}_{jk} - w_{jk}^{*}) \right| \\ &\leq \left\| \hat{w}_{jk} - w_{jk}^{*} \right\|_{1}^{2} \| \Sigma_{j \setminus k, j \setminus k} \|_{\infty} + 2 \| \Sigma_{j \setminus k, j \setminus k} w_{jk}^{*} \|_{\infty} \| \hat{w}_{jk} - w_{jk}^{*} \|_{1} \\ &= \mathcal{O}_{\mathbb{P}} \{ (s_{\max}' \lambda_{w})^{2} \} + \mathcal{O}_{\mathbb{P}} \{ (s_{\max}' \lambda_{w}) w_{0} \} \\ &= \mathcal{O}_{\mathbb{P}} \{ (s_{\max}' \lambda_{w}) w_{0} \}, \end{split}$$

where the last two equalities are obtained by an application of Lemma S9, the fact that  $\|\Sigma\|_{\infty} = O(1)$ , and Assumption S1. Similar upper bounds can be obtained for  $I_{2k}$  and  $I_{3k}$ .

Upon combining the upper bounds, we have

$$\begin{aligned} |\hat{\sigma}_{jk}^{2} - \sigma_{jk}^{2}| &\leq I_{1} + I_{2j} + I_{2k} + I_{3j} + I_{3k} \\ &= \mathcal{O}_{\mathbb{P}}\{(1 + w_{0} + w_{0}^{2})s_{\max}\lambda\log^{6}p/\kappa_{\min}^{2}\} + \mathcal{O}_{\mathbb{P}}(w_{0}s_{\max}'\lambda_{w}) \\ &= o_{\mathbb{P}}(1), \end{aligned}$$

<sup>495</sup> where the last expression holds by Assumption S2.

## S8. TAIL INEQUALITIES

In this section, we state some results that are frequently used in our proofs.

LEMMA S13 (GAUSSIAN TAIL INEQUALITY). Let  $X \sim N(0, \sigma^2)$ . Then, for x > 0,

$$\operatorname{pr}(|X| \ge x) \leqslant \frac{2\sigma}{x} \exp(-x^2/2\sigma^2).$$

LEMMA S14 (LEMMA H.3 IN NING & LIU, 2016). Let  $X_1, \ldots, X_n$  be independent and identically distributed random variables with  $E(X_i) = 0$ . Let  $\overline{X}_n = \sum_{i=1}^n X_i/n$ . If there exist constants  $L_1$ ,  $L_2$  and q such that for x > 0,

$$\operatorname{pr}(|X_i| \ge x) \le L_1 \exp(-L_2 x^q),$$

then for  $x \ge \{8E(X_i^2)/n\}^{1/2}$ ,

$$\operatorname{pr}(|\bar{X}_n| \ge x) \le 4 \exp\left\{-\frac{1}{8}n^{q/(2+q)}x^{2q/(2+q)}\right\} + 4nL_1 \exp\left\{-\frac{L_2 n^{q/(2+q)}x^{2q/(2+q)}}{2^q}\right\} + 4nL_1 \exp\left\{-\frac{L_2 n^{q/(2+q)}x^{2q/(2+q)}}{2^q}\right\} + 4nL_1 \exp\left\{-\frac{L_2 n^{q/(2+q)}x^{2q/(2+q)}}{2^q}\right\} + 4nL_1 \exp\left\{-\frac{L_2 n^{q/(2+q)}x^{2q/(2+q)}}{2^q}\right\} + 4nL_1 \exp\left\{-\frac{L_2 n^{q/(2+q)}x^{2q/(2+q)}}{2^q}\right\}$$

#### REFERENCES

BICKEL, P. J., RITOV, Y. & TSYBAKOV, A. B. (2009). Simultaneous analysis of Lasso and Dantzig selector. Ann. Statist. 37, 1705–32.

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.

510 VAN DE GEER, S. A. & BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the lasso. *Electron. J. Statist.* 3, 1360–92.

YANG, E., RAVIKUMAR, P., ALLEN, G. I. & LIU, Z. (2015). Graphical models via univariate exponential family distributions. J. Mach. Learn. Res. 16, 3813–47.

YE, F. & ZHANG, C.-H. (2010). Rate minimaxity of the lasso and Dantzig selector for the  $\ell_q$  loss in  $\ell_r$  balls. J.

515 *Mach. Learn. Res.* **11**, 3519–40.

505

Latent variable graphical models