# Bayesian joint inference for multiple directed acyclic graphs

Kyoungjae Lee[1] and Xuan Cao[2]

[1]Department of Statistics, Inha university
[2]Department of Mathematical Sciences, University of Cincinnati

August 17, 2020

## Abstract

In many applications, data often arise from multiple groups that may share similar characteristics. A joint estimation method that models several groups simultaneously can be more efficient than estimating parameters in each group separately. We focus on unraveling the dependence structures of data based on directed acyclic graphs and propose a Bayesian joint inference method for multiple graphs. To encourage similar dependence structures across all groups, a Markov random field prior is adopted. We establish the joint selection consistency of the fractional posterior in high dimensions, and benefits of the joint inference are shown under the common support assumption. This is the first Bayesian method for joint estimation of multiple directed acyclic graphs. The performance of the proposed method is demonstrated using simulation studies, and it is shown that our joint inference outperforms other competitors. We apply our method to an fMRI data for simultaneously inferring multiple brain functional networks.

Key words: Joint selection consistency, Markov random field prior, Cholesky factor

## 1 Introduction

Suppose we observe data from the following $K$ groups,

$$X_{k,1}, \ldots, X_{k,n_k} \mid \Omega_k \overset{ind.}{\sim} N_p(0, \Omega_k^{-1}), \ k = 1, \ldots, K, \tag{1}$$

where $\Omega_k \in \mathbb{R}^{p \times p}$ is the precision matrix of the $k$th group. Here, $N_p(\mu, \Sigma)$ denotes the $p$-dimensional normal distribution with the mean vector $\mu \in \mathbb{R}^p$ and covariance matrix $\Sigma \in \mathbb{R}^{p \times p}$. We are interested in investigating the dependence structures of each multivariate data set, especially in high-dimensional settings. To consistently recover the dependence structure of multivariate data, various sparsity assumptions have been suggested for high-dimensional covariance matrices (Cai et al.; 2010; Cai and Zhou; 2012b), precision matrices (Banerjee and Ghosal; 2015; Ren et al.;

2015) and Cholesky factors (Lee and Lee; 2017; Cao et al.; 2019). In this paper, we focus on sparse Cholesky factors, whose sparsity patterns are related to directed acyclic graph (DAG) models. Our goal is to develop a theoretically supported Bayesian method for jointly estimating multiple DAGs under a sparsity assumption.

In many applications, data are collected from multiple groups that share similar characteristics. For examples, gene expression levels are often measured over the patients with different subtypes (Cai et al.; 2016; Liu et al.; 2019), where the DAGs may vary across subtypes but share similar structures. Then, joint estimation can be more efficient than estimating each DAG separately. Another motivation for this type of problem comes from neuroimaging studies. In neuroimaging studies, it is common to explore the changes in functional connectivity for different brain regions through the progression of a certain disease. Taking the Parkinson's disease (PD) as an example, during the progression of PD, some patients may develop the comorbidity of depression, and others may not. Neuroscientists are interested in learning the complex interactions that govern brain connectivity networks and contribute to the onset of depression. In such applications, statistical methods for jointly estimating multiple DAGs can serve as a powerful tool to gain insight into the underlying neurological mechanism.

When data are collected from a homogeneous population, many statistical methods for estimating high-dimensional sparse Cholesky factors have been developed. Shojaie and Michailidis (2010) proposed a penalized likelihood method based on a lasso-type penalty and derived its convergence rate. van de Geer and Bühlmann (2013) showed the convergence rate of the $\ell_0$-penalized maximum likelihood estimator for sparse Cholesky factors. Recently, Khare et al. (2019) developed a convex sparse Cholesky selection, by using a reparameterization trick, and proved the convergence rate and selection consistency in a moderate high-dimensional setting. From a Bayesian perspective, Ben-David et al. (2015) introduced a class of DAG-Wishart priors for sparse DAG models, and Cao et al. (2019) showed the posterior convergence rate and selection consistency of hierarchical DAG-Wishart priors. Based on the autoregressive model representation of a Gaussian DAG model, Lee et al. (2019) developed an empirical sparse Cholesky prior. They showed that the proposed prior attains the minimax optimal posterior convergence rate as well as the selection consistency under mild conditions. However, the above methods are lack of sharing information across graphs when estimating multiple graphs with similar structures.

To infer data sets from heterogeneous populations, various methods have been proposed for estimating multiple graphical models, i.e., precision matrices, by Danaher et al. (2014), Cai et al. (2016), Peterson et al. (2015) and Gan et al. (2019), to name a few. On the other hand, only few joint inference methods for multiple DAGs have been proposed in the literature. Wang et al. (2020) proposed the joint greedy equivalence search for estimating multiple DAGs and proved its convergence rate under the Frobenius norm. They showed that the cardinality of the union of estimated DAGs has the same rate with that of the union of true DAGs. Recently, Liu et al. (2019)

proposed a two-step method, called the multiple PenPC, to jointly estimate the skeletons of DAGs and showed the joint selection consistency of the skeletons in high-dimensional settings. To the best of our knowledge, no Bayesian method, which enjoys theoretical guarantees in high-dimensional settings, has yet been suggested for multiple DAGs.

In this paper, we propose a prior for Bayesian joint inference, called the joint empirical sparse Cholesky prior, for multiple DAGs in high-dimensional settings. We show that the proposed prior achieves the joint selection consistency under mild conditions, which means that the marginal posterior at the true DAGs converges to one as more data are collected (Theorem 3.1). To the best of our knowledge, this is the first work that has established the joint selection consistency for multiple DAGs under a Bayesian framework. We also prove theoretical benefits of the joint inference under the common support assumption. Specifically, it is shown that the proposed method attains the joint selection consistency under much weaker beta-min conditions (Theorems 3.3 and 3.4) compared with separate inferences. In simulation studies, our joint inference method outperforms the other state-of-the-art methods including frequentist joint estimators and Bayesian separate inferences especially in high overlapping scenarios. These finding support our motivation for joint inference: when multiple DAGs share similar structures, joint estimation can be more efficient than separate estimations.

The rest of paper is organized as follows. Section 2 introduces multiple Gaussian DAG models, the joint empirical sparse Cholesky prior and the fractional posterior distribution. In Section 3, we show the joint selection consistency of the proposed method and benefits of the joint inference compared with separate inferences. The finite sample performance of our method is investigated in Section 4, and we conduct a real data analysis using a functional magnetic resonance imaging (fMRI) dataset in Section 5. Section 6 concludes the paper with a discussion. The proofs of the main results are given in Section 7.

## 2 Preliminaries

### 2.1 Multiple Gaussian DAG models

For a given precision matrix $\Omega \in \mathbb{R}^{p \times p}$, let $\Omega = (I_p - A)^T D^{-1}(I_p - A)$ be its modified Cholesky decomposition (MCD), where $A = (a_{jl})$ is a lower triangular matrix with $a_{jj} = 0$ and $D = diag(d_j)$ with $d_j > 0$, for all $j = 1, \ldots, p$. Then, it is well known that $X = (X_1, \ldots, X_p)^T \sim N_p(0, \Omega^{-1})$ can be represented as a sequence of linear autoregressive models as follows:

$$
\begin{aligned}
X_1 \mid d_j &\sim N(0, d_1), \\
X_j \mid a_{S_j}, d_j, S_j &\sim N\Big( \sum_{l \in S_j} X_l a_{jl}, d_j \Big), \ j = 2, \ldots, p,
\end{aligned}
$$

where $a_{S_j} = (a_{jl})_{l \in S_j}^T \in \mathbb{R}^{|S_j|}$, $S_j \subseteq \{1, \ldots, j-1\}$ and $|S_j|$ is the cardinality of $S_j$ (Bickel and Levina; 2008). The support of the Cholesky factor, $\{S_2, \ldots, S_p\}$, determines the DAG, $\mathcal{D} = (V, E)$. Here, $V = \{1, \ldots, p\}$ is a set of vertices, and $E$ is a set of directed edges, where $\{l \to j\} \in E$ if and only if $a_{jl} \neq 0$. In this paper, we assume that a parent ordering of variables is known in which no edges exist from larger vertices to smaller vertices. The above model is called the Gaussian DAG model.

Similarly, for a given $1 \leq k \leq K$, we denote the MCD of $\Omega_k$ by $\Omega_k = (I_p - A_k)^T D_k^{-1}(I_p - A_k)$, where $A_k = (a_{k,jl})$ and $D_k = diag(d_{kj})$. Let $S_{kj} = (S_{k,j1}, \ldots, S_{k,jj-1}) \in \{0, 1\}^{j-1}$ be the support of the $j$th row of $A_k$ with $S_{k,jl} = I(a_{k,jl} \neq 0)$. With a slight abuse of notation, if there is no confusion, $S_{kj}$ is sometimes used to denote the set of nonzero indices in the $j$th row of $A_k$, i.e., $S_{kj} = \{l : a_{k,jl} \neq 0\} \subseteq \{1, \ldots, j-1\}$. We denote the data from the $k$th group and the whole data by $\mathbf{X}_k = (X_{k,1}, \ldots, X_{k,n_k})^T \in \mathbb{R}^{n_k \times p}$ and $\tilde{\mathbf{X}}_n = (\mathbf{X}_1^T, \ldots, \mathbf{X}_K^T)^T \in \mathbb{R}^{n \times p}$, respectively, where $n = \sum_{k=1}^K n_k$. Then, model (1) can be expressed as follows:

$$
\begin{aligned}
\mathbf{X}_{k,1} \mid d_{kj} &\overset{ind.}{\sim} N_{n_k}(0, d_{kj} I_{n_k}), \\
\mathbf{X}_{k,j} \mid a_{k,S_{kj}}, d_{kj}, S_{kj} &\overset{ind.}{\sim} N_{n_k}\Big(\mathbf{X}_{k,S_{kj}} a_{k,S_{kj}}, \; d_{kj} I_{n_k}\Big), \; j = 2, \ldots, p, \; k = 1, \ldots, K,
\end{aligned}
\tag{2}
$$

where $a_{k,S_{kj}} = (a_{k,jl})_{l \in S_{kj}}^T \in \mathbb{R}^{|S_{kj}|}$ and $\mathbf{X}_{k,S} \in \mathbb{R}^{n_k \times |S|}$ is the submatrix consisting of $S$th columns of $\mathbf{X}_k$ for any $S \subseteq \{1, \ldots, j-1\}$. We call model (2) the multiple Gaussian DAG models. Note that the lower triangular part of $A_k$ can be seen as a set of regression vectors, thus we can use a prior tailored to each row of the sparse regression coefficient vectors. We assume that the sample size for each group, $n_k$, can be different across all groups. We consider the high-dimensional setting in which $p \geq n$ and allow the number of groups, $K$, grow to infinity as we observe more data.

## 2.2 Joint empirical sparse Cholesky priors

Lee et al. (2019) proposed the empirical sparse Cholesky (ESC) prior for a sparse DAG model on the basis of the interpretation (2). In this paper, we extend this prior to deal with multiple DAGs. For given $1 \leq k \leq K$ and $S_{kj}$, we use the following conditional prior for $A_k$ and $D_k$:

$$
\begin{aligned}
a_{k,S_{kj}} \mid d_{kj}, S_{kj} &\overset{ind.}{\sim} N_{|S_{kj}|}\Big(\widehat{a}_{k,S_{kj}}, \; \frac{d_{kj}}{\gamma}(\mathbf{X}_{k,S_{kj}}^T \mathbf{X}_{k,S_{kj}})^{-1}\Big), \quad j = 2, \ldots, p, \\
\pi(d_{kj}) &\propto d_{kj}^{-\nu_0/2-1}, \quad j = 1, \ldots, p,
\end{aligned}
\tag{3}
$$

for some positive constants $\gamma$ and $\nu_0$, where $\widehat{a}_{k,S_{kj}} = (\mathbf{X}_{k,S_{kj}}^T \mathbf{X}_{k,S_{kj}})^{-1} \mathbf{X}_{k,S_{kj}}^T \mathbf{X}_{k,j}$. This corresponds to the ESC prior when $K = 1$. Note that the conditional prior for $a_{k,S_{kj}}$ is an empirical version of the Zellner's $g$-prior (Zellner; 1986) centered at $\widehat{a}_{k,S_{kj}}$, and the prior for $d_{kj}$ becomes the Jeffreys prior (Jeffreys; 1946) when $\nu_0 = 0$.

For joint inference on multiple DAGs, given an integer $j \in \{2, \ldots, p\}$, we propose the following joint prior for $(S_{1j}, \ldots, S_{Kj})$:

$$\pi(S_{1j}, \ldots, S_{Kj}) \quad \propto \quad f(S_{1j}, \ldots, S_{Kj}) \prod_{k=1}^{K} \pi(S_{kj}), \tag{4}$$

where

$$\pi(S_{kj}) \quad \propto \quad \binom{j-1}{|S_{kj}|}^{-1} p^{-c_1|S_{kj}|} I(0 \leq |S_{kj}| \leq R_j)$$

for some positive integers $0 < R_j \leq j-1$. Here, $\pi(S_{kj})$ plays a role as a penalty term for the model size $|S_{kj}|$, which prefers sparse models. Similar priors have been used in the literature including Martin et al. (2017) and Lee et al. (2019). For $f(S_{1j}, \ldots, S_{Kj})$ in (4), we suggest using the following Markov random field (MRF) type prior to reflect the expectation that different groups share similar DAG structures:

$$
\begin{aligned}
f(S_{1j}, \ldots, S_{Kj}) \quad &= \quad \exp\left\{ c_{2j} \sum_{l=1}^{j-1} \tilde{S}_{jl}^T (1_K 1_K^T - I_K) \tilde{S}_{jl} \right\} \\
&= \quad \exp\left\{ 2c_{2j} \sum_{l=1}^{j-1} \sum_{k<k'} I(S_{k,jl} = S_{k',jl} = 1) \right\}, \quad j = 2, \ldots, p
\end{aligned}
$$

for some constant $c_{2j} > 0$, where $\tilde{S}_{jl} = (S_{1,jl}, \ldots, S_{K,jl})^T$ and $1_K = (1, \ldots, 1)^T \in \mathbb{R}^K$. This MRF prior encourages similar patterns of sparsity for $(S_{1j}, \ldots, S_{Kj})$. Peterson et al. (2015) used a similar MRF prior for inferring multiple graphical models. By putting together priors (3) and (4), we propose a prior for multiple DAGs,

$$\pi(\Omega_1, \ldots, \Omega_K) \quad \propto \quad \prod_{j=2}^{p} \pi(S_{1j}, \ldots, S_{Kj}) \prod_{k=1}^{K} \left\{ \prod_{j=2}^{p} \pi(a_{k,S_{kj}} \mid d_{kj}, S_{kj}) \prod_{j=1}^{p} \pi(d_{kj}) \right\},$$

which we call the joint empirical sparse Cholesky (JESC) prior hearafter.

## 2.3   $\alpha$-fractional posterior

We adopt the fractional likelihood framework, which has received increasing attention in recent years (Martin and Walker; 2014; Martin et al.; 2017; Lee et al.; 2019). Let $\theta$ and $L(\theta)$ be a parameter and a likelihood function, respectively. For a given constant $\alpha \in (0,1)$, $\alpha$-fractional likelihood $L_\alpha(\theta)$ is the likelihood with power $\alpha$, i.e., $\{L(\theta)\}^\alpha$. Based on the JESC prior and $\alpha$-fractional likelihood, we have the following posterior distributions:

$$a_{k,S_{kj}} \mid d_{kj}, S_{kj}, \mathbf{X}_k \overset{ind.}{\sim} N_{|S_{kj}|}\left( \widehat{a}_{k,S_{kj}}, \frac{d_{kj}}{\alpha + \gamma} (\mathbf{X}_{k,S_{kj}}^T \mathbf{X}_{k,S_{kj}})^{-1} \right), \ j = 2, \ldots, p,$$

$$d_{kj} \mid S_{kj}, \mathbf{X}_k \overset{ind.}{\sim} IG\left( \frac{\alpha n_k + \nu_0}{2}, \frac{\alpha n_k}{2} \widehat{d}_{k,S_{kj}} \right), \ j = 1, \ldots, p,$$

5

and

$$\pi_\alpha(S_{1j}, \ldots, S_{Kj} \mid \tilde{\mathbf{X}}_n) \quad \propto \quad \pi(S_{1j}, \ldots, S_{Kj}) \prod_{k=1}^K f_\alpha(\mathbf{X}_{n_k} \mid S_{kj}), \ j = 2, \ldots, p,$$

where $\hat{d}_{k,S_{kj}} = n_k^{-1} \mathbf{X}_{k,j}^T (I_{n_k} - \tilde{P}_{S_{kj}}) \mathbf{X}_{k,j}$, $\tilde{P}_{S_{kj}} = \mathbf{X}_{k,S_{kj}} (\mathbf{X}_{k,S_{kj}}^T \mathbf{X}_{k,S_{kj}})^{-1} \mathbf{X}_{k,S_{kj}}^T$ and

$$\begin{aligned}
f_\alpha(\mathbf{X}_{n_k} \mid S_{kj}) &= \iint L_\alpha(a_{k,S_{kj}}, d_{kj}, S_{kj}) \pi(a_{k,S_{kj}} \mid d_{kj}, S_{kj}) \pi(d_{kj} \mid S_{kj}) da_{k,S_{kj}} \, dd_{kj} \\
&\propto \left(1 + \frac{\alpha}{\gamma}\right)^{-\frac{|S_{kj}|}{2}} (\hat{d}_{k,S_{kj}})^{-\frac{\alpha n_k + \nu_0}{2}}.
\end{aligned}$$

We denote the posterior by $\pi_\alpha(\cdot \mid \tilde{\mathbf{X}}_n)$ to indicate that the $\alpha$-fractional likelihood is used, and call it the $\alpha$-fractional posterior. To conduct the posterior inference for $(S_{1j}, \ldots, S_{Kj})$, the Metropolis-Hastings within Gibbs algorithm can be used. The details are given in Section 4.1. Once we have posterior samples of $(S_{1j}, \ldots, S_{Kj})$, the posterior samples of $a_{k,S_{kj}}$ and $d_{kj}$ can be directly drawn from the normal and inverse-gamma distributions, respectively.

# 3 Main Results

## 3.1 Joint selection consistency

In this section, we establish the joint selection consistency of the proposed JESC prior, which guarantees that we can recover the true DAGs asymptotically. Let $\Omega_{0k}$ be the true precision matrix of the $k$th class, for $k = 1, \ldots, K$. Let $\Omega_{0k} = (I_p - A_{0k})^T D_{0k}^{-1}(I_p - A_{0k})$ be the MCD of $\Omega_{0k}$, where $A_{0k} = (a_{0k,jl})$ and $D_{0k} = diag(d_{0k,j})$. We denote $S_A$ as the support of the matrix $A = (a_{jl})$, i.e., $S_A = (I(a_{jl} \neq 0))$. We first introduce the following sufficient conditions for true parameters:

**Condition (A1)** There exists a constant $0 < \epsilon_0 < 0.5$ such that $\epsilon_0 \leq \min_{1 \leq k \leq K} \lambda_{\min}(\Omega_{0k}) \leq \max_{1 \leq k \leq K} \lambda_{\max}(\Omega_{0k}) \leq \epsilon_0^{-1}$.

**Condition (A2)** $\max_{1 \leq k \leq K} \max_{2 \leq j \leq p} \sum_{l=1}^p I(a_{0k,jl} \neq 0) \leq s_0$ for some $1 \leq s_0 \leq p$.

**Condition (A3)** For some constant $C_{\mathrm{bm}} > 0$,

$$\min_{1 \leq k \leq K} \min_{(j,l):a_{0k,jl} \neq 0} n_k a_{0k,jl}^2 \quad \geq \quad \frac{16}{\alpha(1-\alpha)\epsilon_0^2(1-2\epsilon_0)^2} C_{\mathrm{bm}} \log p.$$

**Condition (A4)** $K = o(\log p)$.

Condition (A1) implies that the eigenvalues of each precision matrix $\Omega_{0k}$ are bounded. This condition is used to obtain upper bounds of $d_{0k,j}, d_{0k,j}^{-1}$ and $\|A_{0k}\|$. Similar conditions have been used in, for examples, Ren et al. (2015), Khare et al. (2019) and Lee et al. (2019).

Condition (A2) controls the maximum number of nonzero entries in each row of $A_{0k}$. This condition allows the upper bound $s_0$ to grow to infinity as $n$ get larger. Note that the estimation of each row of $A_{0k}$ can be considered as the estimation of regression coefficient vector, thus introducing this condition seems natural.

Condition (A3) is the well known *beta-min* condition for the minimum nonzero entries of each Cholesky factor, $A_{0k}$. This roughly means that the lower bound for nonzero $a_{0k,jl}^2$ is of order $O(\log p/n_k)$. The beta-min condition is essential for consistent variable selection in high-dimensional linear regression models (Martin et al.; 2017; Yang et al.; 2016) and Gaussian DAG models (Yu and Bien; 2017; Cao et al.; 2019). Note that if we assume $k = 1$ and $n_k = n$, then the rate of the lower bound in condition (A3) becomes $\log p/n$, which is the best (minimum) beta-min condition in the literature.

Condition (A4) restricts the number of classes. Note that $K$ can grow to infinity as $n \to \infty$ at a rate slower than $\log p$. Cai et al. (2016) and Wang et al. (2020) used similar condition for joint estimation of high-dimensional precision matrices and DAGs, respectively.

**Condition (P)** $\nu_0 = o(\min_k n_k), c_1 > 2, c_{2j} \leq 1/(j-1)$ and $\gamma = O(1)$. For some small $0 < c_3 < (\epsilon')^2 \epsilon_0^2/\{128(1+2\epsilon_0)^2\}$ and $\epsilon' = \{(1-\alpha)/10\}^2$, we assume that $R_j = \lfloor \{(\log n)^{-1} \vee c_3\} \min_k n_k/\log p \rfloor$.

Condition (P) shows a sufficient condition for hyperparameters in the JESC prior to obtain the desired theoretical properties, where "P" stands for "prior". The constant $c_1$ controls the penalty for the sparsity of Cholesky factors, thus the condition $c_1 > 2$ gives the minimum strength of the penalty. The constant $c_{2j}$ in the MRF prior controls the penalty for similarities across the DAGs, thus $c_{2j} \leq 1/(j-1)$ implies that the effect of the MRF prior should not be too strong. This intuitively makes sense because if $c_{2j}$ is too large and dominates the other priors and likelihoods, then the posterior will always select the full model, i.e., $S_{kj} = \{1, \ldots, j-1\}$ for all $1 \leq k \leq K$ and $2 \leq j \leq p$. The condition $R_j = \lfloor \{(\log n)^{-1} \vee c_3\} \min_k n_k/\log p \rfloor$ implies that the maximum number of nonzero entries in each row of $A_{0k}$ should at least be of order $\min_k n_k/\log p$ for the consistent selection. In finite samples, we suggest choosing $R_j = \lfloor \min_k n_k(\log p \log n)^{-1} \rfloor$. In Section 4.1, we will give a practical guidance for the choice of hyperparameters.

**Theorem 3.1 (Joint selection consistency)** *Suppose that conditions (A1)-(A4) and (P) hold with $C_{\mathrm{bm}} > c_1 + 2$. Then, if $s_0 \log p \leq \min_k n_k c_3/2$ and $s_0 \geq C_{\mathrm{bm}} - c_1 - 1$, we have*

$$\mathbb{E}_0\Big\{\pi_\alpha\Big(S_{A_1} = S_{A_{01}}, \ldots, S_{A_K} = S_{A_{0K}} \mid \tilde{\mathbf{X}}_n\Big)\Big\} \longrightarrow 1 \quad as \ \min_k n_k \to \infty.$$

Theorem 3.1 presents the joint selection consistency for multiple DAGs. It is worth comparing our result with those in Liu et al. (2019) in terms of the required conditions. To obtain consistency, they assumed $\min_k \lambda_{\min}(\Sigma_{0k,AA}) \geq C_1$ and $\max_k \Sigma_{0k,jj} < C_2$ for any $A \in \{1, \ldots, p\}$ with $|A| \leq$

$q$ and some constants $C_1$ and $C_2 > 0$, which is weaker than condition (A1), where $\Sigma_{0k,AA} = (\Sigma_{0k,jl})_{j,l \in A}$. It was also assumed that $n_k \asymp n$ for all $k = 1, \ldots, K$. Note that this condition implies $K = O(1)$, thus it is stronger than our condition (A4). They further assumed $p = O(\exp(n^a))$ and $q = \max_j |A_j| = O(n^b)$ for some constants $a \in [0, 1)$ and $b \in [0, (1-a)/2)$, where $A_j = \cup_{k=1}^{K} A_j^{(k)}$ and $A_j^{(k)} = \{l : \Omega_{0k,jl} \neq 0 \text{ and } l \neq j\}$. By Lemma 1 in Liu et al. (2019), $q = O(n^b)$ implies $s_0 \leq |\cup_{k=1}^{K} S_{0k,j}| = O(n^b)$, thus it is slightly more restrictive than our conditions, (A2) and $s_0 \log p \leq \min_k n_k c_3 / 2$. They used the beta-min conditions,

$$\min_{1 \leq k \leq K} \min_{(j,l):\Omega_{0k,jl} \neq 0} \left| \frac{\Omega_{0k,jl}}{\Omega_{0k,jj}} \right| \gtrsim n^{-d_1} \quad \text{and} \quad n^{-d_2} \lesssim |\rho_{jl|S}^{(k)}| \leq M < 1, \tag{5}$$

for some $0 < d_1 < (1 - a - b)/2$, $0 < d_2 < \{1 - (a \vee b)\}/2$ and any $S \in \Pi_{jl}^{(k)}$, where $\Pi_{jl}^{(k)} = \{A_{jl}^{(k)} \setminus D_{jl}^{(k)} : D_{jl}^{(k)} \subseteq C_{jl}^{(k)}\}$, $A_{jl}^{(k)}$ is the Markov blanket of $j$ and $l$ after removing their common children and descendants, and $C_{jl}^{(k)}$ is the set of common children or descendants. Note that (5) consists of two beta-min conditions to guarantee selection consistency in each step. Although their beta-min conditions are not directly comparable with ours, the squares of the lower bounds in (5), $n^{-2d_1}$ and $n^{-2d_2}$, are much larger than $\log p / n_k$ in condition (A3). Therefore, we obtain the joint selection consistency under weaker conditions on $K$, $s_0$ and minimum nonzero signals than those in Liu et al. (2019).

**Theorem 3.2** *Let* $\pi^I(S_{kj} \mid \mathbf{X}_{n_k}) \propto f_\alpha(\mathbf{X}_{n_k} \mid S_{kj})\pi(S_{kj})$ *be the independence posterior for* $S_{kj}$. *Suppose that there exists* $1 \leq k \leq K$ *such that* $\cup_{k' \neq k} S_{0k',j} \subseteq S_{0k,j}$. *Then, for any* $j = 2 \ldots, p$, *we have*

$$\pi_\alpha\big(S_{0k,j} \mid S_{01,j}, \ldots, S_{0k-1,j}, S_{0k+1,j}, \ldots, S_{0K,j}, \tilde{\mathbf{X}}_n\big) \geq \pi_\alpha^I(S_{0k,j} \mid \mathbf{X}_{n_k}). \tag{6}$$

Theorem 3.2 shows that the joint inference increases the conditional posterior probability at the true DAGs compared to the separate inference. Note that $\cup_{k' \neq k} S_{0k',j} \subseteq S_{0k,j}$ holds if and only if

$$f(S_{01,j}, \ldots, S_{0k-1,j}, S_{kj}, S_{0k+1,j}, \ldots, S_{0K,j}) \leq f(S_{01,j}, \ldots, S_{0k-1,j}, S_{0k,j}, S_{0k+1,j}, \ldots, S_{0K,j}) \tag{7}$$

for any $S_{kj} \neq S_{0k,j}$. For example, (7) trivially holds if we assume the common support, i.e., $S_{01,j} = \cdots = S_{0K,j}$.

## 3.2 Benefits of joint inference

In this section, the theoretical benefits of the joint inference, compared with separate inferences, are presented. Although investigating benefits of the joint inference under heterogeneous DAGs is important, it is very challenging to explore every possible scenario. Thus, we focus on the case where all Cholesky factors share a common support, i.e., all DAGs share a common structure. For example, Cai et al. (2016) and Gan et al. (2019) also used the common support assumption for

multiple precision matrices and showed advantages of the joint estimation. In this case, we suggest using the restricted posterior to the space of common supports,

$$\tilde{\pi}_\alpha(S_A \mid \tilde{\mathbf{X}}_n) = \frac{\pi_\alpha(S_{A_1} = \cdots = S_{A_K} = S_A \mid \tilde{\mathbf{X}}_n)}{\sum_{S_A} \pi_\alpha(S_{A_1} = \cdots = S_{A_K} = S_A \mid \tilde{\mathbf{X}}_n)}.$$

To prove the joint selection consistency of $\tilde{\pi}_\alpha(S_A \mid \tilde{\mathbf{X}}_n)$, we introduce a weakened beta-min condition as follows:

**Condition (B3)** For some constant $C_{\mathrm{bm}} > 0$,

$$\min_{(j,l):a_{01,jl}\neq 0} \sum_{k=1}^{K} n_k\, a_{0k,jl}^2 \geq \frac{16}{\alpha(1-\alpha)\epsilon_0^2(1-2\epsilon_0)^2} C_{\mathrm{bm}} K \log p.$$

Note that condition (A3) implies (B3), thus we call condition (B3) a weakened beta-min condition. If we assume that $n_1 = \cdots = n_K$, then condition (B3) roughly means that the lower bound for $K^{-1}\sum_{k=1}^{K} \min_{(j,l):a_{0k,jl}\neq 0} a_{0k,jl}^2$ is of order $O(\log p/n_k)$. Thus, we can consistently recover the true support as long as the *average* of minimum signals is significant, even if minimum signals of some classes are quite small. This can be seen as the benefit of the joint inference, and the following theorem states the desired result.

**Theorem 3.3 (Benefit of joint inference)** *Assume that $S_{A_{01}} = \cdots = S_{A_{0K}} \equiv S_0$ and $K\log p = o(\min_k n_k)$. Then, under the same condition with Theorem 3.1, except using condition (B3) instead of (A3), we have*

$$\mathbb{E}_0\big\{\tilde{\pi}_\alpha(S_A = S_0 \mid \tilde{\mathbf{X}}_n)\big\} \longrightarrow 1 \quad \text{as } \min_k n_k \to \infty.$$

Note that $K\log p = o(\min_k n_k)$ trivially holds if we assume $(\log p)^2 = o(\min_k n_k)$, by condition (A4). Cai et al. (2016) assumed $K^{2a-1}\log p\,(\log n)^2 = o(\min_k n_k)$ and $\max(K, K^{4-a}\log K) = o(\log p)$ for some constant $a > 0$. The second condition is comparable to our condition (A4) when $a = 3$, and then the first condition becomes $K^5 \log p\,(\log n)^2 = o(\min_k n_k)$. Thus, our condition $K\log p = o(\min_k n_k)$ is much weaker than that of Cai et al. (2016).

In fact, if we slightly modify the prior for $S_A$, we can further weaken the beta-min condition. Define the modified prior for $(S_{1j}, \ldots, S_{Kj})$ as

$$\begin{aligned}
\tilde{\pi}(S_{1j}, \ldots, S_{Kj}) &\propto \pi(S_{1j}, \ldots, S_{Kj})^{1/K} \\
&\propto f(S_{1j}, \ldots, S_{Kj})^{1/K} \prod_{k=1}^{K} \pi(S_{kj})^{1/K} \\
&\equiv \tilde{f}(S_{1j}, \ldots, S_{Kj}) \prod_{k=1}^{K} \tilde{\pi}(S_{kj}),
\end{aligned}$$

9

and let $\tilde{\pi}_\alpha^*(S_A \mid \tilde{\mathbf{X}}_n) = \tilde{\pi}_\alpha^*(S_{A_1} = \cdots = S_{A_K} = S_A \mid \tilde{\mathbf{X}}_n)$ be the restricted posterior to the space of common supports using the prior $\tilde{\pi}(S_{1j}, \ldots, S_{Kj})$ instead of $\pi(S_{1j}, \ldots, S_{Kj})$. Then, it suffices to assume the following condition (C3) instead of condition (B3) to obtain the joint selection consistency:

**Condition (C3)** For some constant $C_{\mathrm{bm}} > 0$,

$$
\min_{(j,l):a_{01,jl} \neq 0} \sum_{k=1}^{K} n_k \, a_{0k,jl}^2 \quad \geq \quad \frac{16}{\alpha(1-\alpha)\epsilon_0^2(1-2\epsilon_0)^2} C_{\mathrm{bm}} \log p.
$$

**Theorem 3.4 (Benefit of joint inference II)** *Assume that $S_{A_{01}} = \cdots = S_{A_{0K}} \equiv S_0$. Then, under the same condition with Theorem 3.1, except using condition (C3) instead of (A3), we have*

$$
\mathbb{E}_0\big\{\tilde{\pi}_\alpha^*(S_A = S_0 \mid \tilde{\mathbf{X}}_n)\big\} \quad \longrightarrow \quad 1 \qquad as \ \min_k n_k \to \infty.
$$

Theorem 3.4 shows the advantage of the joint inference based on the restricted posterior $\tilde{\pi}_\alpha^*(S_A \mid \tilde{\mathbf{X}}_n)$: it only requires condition (C3), which is much weaker than condition (B3). Compared with Theorems 3.1 and 3.3, it reveals that, under the common support assumption, we can obtain the joint selection consistency as long as the *summation* of minimun signals is significant. Note that the lower bound in condition (C3) coincides with that in (A3). Cai et al. (2016) used a similar beta-min condition to condition (C3) for the nonzero entries of precision matrices, but using $\log K \log p$ instead of $\log p$. Hence, our beta-min condition is weaker than their in terms of the rate. Also note that $\prod_{k=1}^{K} \tilde{\pi}(S_{kj}) \propto \pi(S_{1j}) I(S_{1j} = \cdots = S_{Kj})$ when $S_{1j} = \cdots = S_{Kj}$. Thus, this implies that it is sufficient to use a single penalty (prior) for all $K$ classes rather than use a penalty for each class.

# 4 Simulation Studies

In this section, we carry out simulation studies to illustrate the model selection performance of our method and show its potential benefits over other contenders.

## 4.1 Posterior inference

The use of the JESC prior not only guarantees the asymptotic properties but also allows us to easily conduct the posterior inference. Recall that for $j = 2, \ldots, p$,

$$
\begin{aligned}
&\pi_\alpha(S_{1j}, \ldots, S_{Kj} \mid \tilde{\mathbf{X}}_n) \\
\propto \ & \prod_{k=1}^{K} \left(1 + \frac{\alpha}{\gamma}\right)^{-\frac{|S_{kj}|}{2}} \big(\hat{d}_{k,S_{kj}}\big)^{-\frac{\alpha n_k + \nu_0}{2}} \binom{j-1}{|S_{kj}|}^{-1} p^{-c_1|S_{kj}|} I\big(0 \leq |S_{kj}| \leq R_j\big) \\
& \times \exp\left\{c_2 \sum_{l=1}^{j-1} \tilde{S}_{jl}^T (1_K 1_K^T - I_K) \tilde{S}_{jl}\right\},
\end{aligned}
$$

where $\tilde{S}_{jl} = (S_{1,jl}, \ldots, S_{K,jl})^T$. Hence, we can run the Metropolis-Hastings within Gibbs sampling algorithm for each $j = 2, \ldots, p$ in parallel. Here, we briefly summarize the algorithm used for the inference:

Run the following steps for $j = 2, \ldots, p$.

1. Set the initial values $S_{1j}^{(1)}, \ldots, S_{Kj}^{(1)}$.

2. For each $t = 2, \ldots, T$, run the following steps for $k = 1, \ldots, K$.

   (a) sample $S_{kj}^{new} \sim q(\cdot \mid S_{kj}^{(t)})$;

   (b) set $S_{kj}^{(t)} = S_{kj}^{new}$ with the probability

   $$
   \min\left\{1, \frac{\pi_\alpha(S_{kj}^{new} \mid S_{1j}^{(t)}, \ldots, S_{k-1,j}^{(t)}, S_{k+1,j}^{(t-1)}, \ldots, S_{Kj}^{(t-1)}, \tilde{\mathbf{X}}_n)q(S_{kj}^{(j-1)} \mid S_{kj}^{new})}{\pi_\alpha(S_{kj}^{(t-1)} \mid S_{1j}^{(t)}, \ldots, S_{k-1,j}^{(t)}, S_{k+1,j}^{(t-1)}, \ldots, S_{Kj}^{(t-1)}, \tilde{\mathbf{X}}_n)q(S_{kj}^{new} \mid S_{kj}^{(t-1)})}\right\}
   $$
   $$
   = \min\left\{1, \frac{\pi_\alpha(S_{kj}^{new} \mid \mathbf{X}_{n_k})f(S_{1j}^{(t)}, \ldots, S_{k-1,j}^{(t)}, S_{kj}^{new}, S_{k+1,j}^{(t-1)}, \ldots, S_{Kj}^{(t-1)})q(S_{kj}^{(j-1)} \mid S_{kj}^{new})}{\pi_\alpha(S_{kj}^{(t-1)} \mid \mathbf{X}_{n_k})f(S_{1j}^{(t)}, \ldots, S_{k-1,j}^{(t)}, S_{kj}^{(t-1)}, S_{k+1,j}^{(t-1)}, \ldots, S_{Kj}^{(t-1)})q(S_{kj}^{new} \mid S_{kj}^{(t-1)})}\right\},
   $$

   otherwise set $S_{kj}^{(t)} = S_{kj}^{(t-1)}$.

The kernel $q(S^{new} \mid S)$ is chosen to form a new set $S^{new}$ by changing a randomly selected nonzero component to 0 with probability 0.5 or by changing a randomly selected zero component to 1 with probability 0.5. Steps 1 and 2 in the above algorithm, can be parallelized for each column. For more details, we refer the interested readers to Cao et al. (2019) and Lee et al. (2019).

The tuning parameters are chosen as suggested in Martin et al. (2017) and Lee et al. (2019). Specifically, we set $\alpha = 0.999$ to mimic the Bayesian model with the original likelihood. In practice, as long as $1 - \alpha$ is close to zero, the performance was not sensitive to the choice of $\alpha$. The other hyperparameters were chosen as $\gamma = 0.1$, $\nu_0 = 0$, $c_1 = 2$ and $c_{2j} = \{p(K-1)\}^{-1}$ for $j = 2, \ldots, p$ to satisfy the theoretical conditions. The above algorithm is coded in R and publicly available at `https://github.com/xuan-cao/Multiple-DAG-Selection`.

## 4.2 Simulation setting

In this section, we demonstrate the performance of the proposed method in various settings similar to those used in Liu et al. (2019); Peterson et al. (2015, 2020). We construct three Cholesky factors $A_1$, $A_2$, and $A_3$ corresponding to DAGs $\mathcal{D}_1$, $\mathcal{D}_2$ and $\mathcal{D}_3$ with different degrees of shared structure. We include $p = 150$ nodes, and consider the first scenario as follows. For the first $p \times p$ lower triangular matrix $A_1$, we randomly chose 2% of the lower triangular entries of $A_1$ and sampled their values from a uniform distribution on $[-0.7, -0.3] \cup [0.3, 0.7]$. The remaining entries were set to zero. $\mathcal{D}_1$ can be acquired by mapping the nonzero entries in $A_1$ to a DAG with $p$ nodes. To

obtain $\mathcal{D}_2$, five edges are removed from $\mathcal{D}_1$ and five new edges added at random. To obtain $\mathcal{D}_3$, five edges are removed from the graph for group 2, and five edges added at random. All the lower triangular entries in $A_2$ and $A_3$ are generated in a similar manner as in $A_1$. We call this simulation setting Scenario 1 (*high overlapping*), where each pair of DAGs have 218 of 223 edges (97.76%) in common.

Next, we investigate a different simulation scenario, say Scenario 2 (*medium overlapping*), where $A_1, \mathcal{D}_1, A_2, \mathcal{D}_2$ are formed as in Scenario 1, but we change the design of $A_3$ and $\mathcal{D}_3$ as follows. To obtain $\mathcal{D}_3$, 20 edges are removed from the graph for group 2, and 20 edges added at random. All the entries in three Cholesky factors $A_1$, $A_2$, and $A_3$ are generated as in Scenario 1. Under this setting, $\mathcal{D}_1$ and $\mathcal{D}_2$ share 218 of 223 edges (97.76%), $\mathcal{D}_2$ and $\mathcal{D}_3$ share 203 edges (91.03%), and $\mathcal{D}_1$ and $\mathcal{D}_3$ share around 219 edges (89.24%). For our final simulation setting, Scenario 3 (*low overlapping*), we first create $A_1$ and $\mathcal{D}_1$ as previously mentioned, and obtain $\mathcal{D}_2$ by randomly removing 20 edges and adding 20 edges from $\mathcal{D}_1$. $\mathcal{D}_3$ is again acquired by randomly removing 20 edges and adding 20 edges from $\mathcal{D}_2$. These steps result in DAGs $\mathcal{D}_1$ and $\mathcal{D}_2$ that share 203 of 223 edges (91.03%), $\mathcal{D}_2$ and $\mathcal{D}_3$ that share 203 edges (91.03%), and $\mathcal{D}_1$ and $\mathcal{D}_3$ that have 185 common edges (82.96%). All the nonzero entries in $A_1$, $A_2$, and $A_3$ are then sampled from a uniform distribution as elaborated in Scenario 1. For all settings, we simulate the diagonal entries of $D_1, D_2, D_3$ from a uniform distribution on $[2, 5]$. Given the precision matrices $\Omega_k = (I_p - A_k)^T D_k^{-1}(I_p - A_k)$ for $k = 1, 2, 3$, the data sets were generated from the multivariate normal distribution $N_p(0, \Omega_k^{-1})$ with $(n_k, p) = (100, 150)$ for $k = 1, 2, 3$.

## 4.3 Performance comparison

We compare the following methods: the proposed JESC prior, Bayesian inference based on ESC applied separately for each group (SESC) (Lee et al.; 2019), multiple PenPC (MPenPC) (Liu et al.; 2019), joint graphical lasso (JGL) (Danaher et al.; 2014), and seperate DAG lasso (DAGL) for each group (Shojaie and Michailidis; 2010). The tuning parameters in JGL were selected using a grid search to identify the combination that minimizes the AIC as suggested in Danaher et al. (2014). Since for our simulation studies, JGL could not produce exact zeros in the Cholesky factors of the estimated precision matrices, we further adopt the hard thresholding of these Cholesky factors. The penalty parameters in MPenPC were tuned using the extended BIC (EBIC) (Chen and Chen; 2008) as suggested in Liu et al. (2019). The penalty parameters in DAGL were set as $\lambda_i(\alpha) = 2n^{-1/2}Z^*_{0.1/\{2p(i-1)\}}$ (separate for each variable $i$), where $Z^*_q$ denotes the $(1-q)^{th}$ quantile of the standard normal distribution. This choice is justified in Shojaie and Michailidis (2010) based on asymptotic considerations. For Bayesian methods, we ran the Metropolis-Hastings algorithm specified in Section 4.1 for each data set to conduct posterior inferences. Every MCMC chain started from an empty initial state and ran for 5,000 iterations with a burn-in period of 1,000, since

Table 1: Performance summary for Scenario 1 (high overlapping). Comparison of true positive rate (TPR), false positive rate (FPR), Matthews correlation coefficient (MCC) and area under the ROC curve (AUC). The models compared are the Bayesian joint ESC method proposed in this paper (JESC) (Lee et al.; 2019), separate ESC method applied for individual group (SESC), multiple PenPC (MPenPC) (Liu et al.; 2019), and joint graphical lasso (JGL) (Danaher et al.; 2014).

|  | Measure | JESC | SESC | MPenPC | JGL | DAGL |
|---|---|---|---|---|---|---|
| Group 1 | TPR | 0.8879 | 0.8610 | 0.8924 | 0.9148 | 0.3785 |
|  | FPR | 0.0045 | 0.0048 | 0.0232 | 0.0365 | 0 |
|  | MCC | 0.8403 | 0.8193 | 0.6163 | 0.5432 | 0.6113 |
|  | AUC | 0.9761 | 0.9684 | . | . | . |
| Group 2 | TPR | 0.9148 | 0.9072 | 0.9462 | 0.9372 | 0.3668 |
|  | FPR | 0.0039 | 0.0044 | 0.0211 | 0.0326 | 0 |
|  | MCC | 0.8664 | 0.8369 | 0.6638 | 0.5769 | 0.6009 |
|  | AUC | 0.9962 | 0.9780 | . | . | . |
| Group 3 | TPR | 0.8969 | 0.8654 | 0.8879 | 0.8789 | 0.3552 |
|  | FPR | 0.0038 | 0.0041 | 0.0230 | 0.0369 | 0 |
|  | MCC | 0.8580 | 0.8361 | 0.6152 | 0.5224 | 0.5901 |
|  | AUC | 0.9835 | 0.9805 | . | . | . |
| All edges | TPR | 0.8999 | 0.8775 | 0.9088 | 0.9103 | 0.3669 |
|  | FPR | 0.0040 | 0.0044 | 0.0224 | 0.0353 | 0 |
|  | MCC | 0.8549 | 0.8308 | 0.6317 | 0.5471 | 0.6010 |
|  | AUC | 0.9479 | 0.9289 | . | . | . |
| Differential edges | TPR | 1 | 0.9050 | 1 | 1 | 0.4600 |
|  | FPR | 0 | 0 | 0 | 0 | 0 |
|  | MCC | 1 | 0.9098 | 1 | 1 | 0.5463 |
|  | AUC | 1 | 0.9525 | . | . | . |

we observed that on average the posterior samples converged rapidly and stabilized after 1,000 iterations. The hyperparameter $c_2$ was set to 0 when implementing SESC. We constructed the final model by collecting indices with inclusion probabilities exceeding 0.5.

To evaluate the performance of joint DAG selection, the true positive rate (TPR), false positive rate (FPR), Matthews correlation coefficient (MCC), and area under the curve (AUC) are reported at Tables 1, 2 and 3 averaged over 20 repetitions. The criteria are defined as

$$
\begin{aligned}
\text{TPR} &= \frac{TP}{TP + FN}, \\
\text{FPR} &= \frac{FP}{TN + FP}, \\
\text{MCC} &= \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}},
\end{aligned}
$$

where TP, TN, FP and FN are true positive, true negative, false positive and false negative, respectively. The AUC is calculated based on the TPR and the FPR for Bayesian methods with varying thresholds. The AUCs for the regularization methods are omitted.

Based on the simulation results (Tables 1 to 3), we can tell that the proposed method is more

Table 2: Performance summary for Scenario 2 (medium overlapping).

|  | Measure | JESC | SESC | MPenPC | JGL | DAGL |
|---|---|---|---|---|---|---|
| Group 1 | TPR | 0.8704 | 0.8475 | 0.9058 | 0.9193 | 0.3565 |
|  | FPR | 0.0051 | 0.0049 | 0.0227 | 0.0363 | 0 |
|  | MCC | 0.8195 | 0.8096 | 0.6275 | 0.5465 | 0.5908 |
|  | AUC | 0.9810 | 0.9836 | · | · | · |
| Group 2 | TPR | 0.9238 | 0.8654 | 0.9238 | 0.9372 | 0.3632 |
|  | FPR | 0.0030 | 0.0043 | 0.0216 | 0.0330 | 0 |
|  | MCC | 0.8901 | 0.8307 | 0.6466 | 0.5748 | 0.5963 |
|  | AUC | 0.9896 | 0.9885 | · | · | · |
| Group 3 | TPR | 0.8610 | 0.8834 | 0.8924 | 0.8879 | 0.3529 |
|  | FPR | 0.0040 | 0.0046 | 0.0232 | 0.0368 | 0 |
|  | MCC | 0.8335 | 0.8359 | 0.6163 | 0.5276 | 0.5897 |
|  | AUC | 0.9859 | 0.9853 | · | · | · |
| All edges | TPR | 0.8849 | 0.8654 | 0.9073 | 0.9148 | 0.3584 |
|  | FPR | 0.0041 | 0.0046 | 0.0225 | 0.0354 | 0 |
|  | MCC | 0.8475 | 0.8254 | 0.6301 | 0.5484 | 0.5930 |
|  | AUC | 0.9405 | 0.9304 | · | · | · |
| Differential edges | TPR | 0.8920 | 0.8482 | 0.9241 | 0.9190 | 0.3800 |
|  | FPR | 0 | 0 | 0.0381 | 0.0814 | 0 |
|  | MCC | 0.8978 | 0.8584 | 0.8867 | 0.8423 | 0.4835 |
|  | AUC | 0.9461 | 0.9235 | · | · | · |

Table 3: Performance summary for Scenario 3 (low overlapping).

|  | Measure | JESC | SESC | MPenPC | JGL | DAGL |
|---|---|---|---|---|---|---|
| Group 1 | TPR | 0.8879 | 0.8520 | 0.8924 | 0.9193 | 0.3796 |
|  | FPR | 0.0042 | 0.0048 | 0.0226 | 0.0360 | 0 |
|  | MCC | 0.8456 | 0.8140 | 0.6207 | 0.5484 | 0.6067 |
|  | AUC | 0.9866 | 0.9786 | · | · | · |
| Group 2 | TPR | 0.8969 | 0.8565 | 0.9148 | 0.9193 | 0.3330 |
|  | FPR | 0.0041 | 0.0040 | 0.0236 | 0.0372 | 0 |
|  | MCC | 0.8526 | 0.8327 | 0.6254 | 0.5422 | 0.5705 |
|  | AUC | 0.9829 | 0.9785 | · | · | · |
| Group 3 | TPR | 0.8879 | 0.9103 | 0.9148 | 0.9148 | 0.3643 |
|  | FPR | 0.0044 | 0.0047 | 0.0255 | 0.0365 | 0 |
|  | MCC | 0.8403 | 0.8497 | 0.6116 | 0.5432 | 0.5967 |
|  | AUC | 0.9869 | 0.985 | · | · | · |
| All edges | TPR | 0.8909 | 0.8729 | 0.9073 | 0.9178 | 0.3576 |
|  | FPR | 0.0042 | 0.0045 | 0.0239 | 0.0366 | 0 |
|  | MCC | 0.8462 | 0.8321 | 0.6191 | 0.5446 | 0.5916 |
|  | AUC | 0.9433 | 0.9342 | · | · | · |
| Differential edges | TPR | 0.8530 | 0.8105 | 0.9255 | 0.8940 | 0.3375 |
|  | FPR | 0 | 0 | 0.0518 | 0.0905 | 0 |
|  | MCC | 0.8617 | 0.8256 | 0.8753 | 0.8392 | 0.4459 |
|  | AUC | 0.9251 | 0.9022 | · | · | · |

conservative in the identification of differential edges compared with frenquentist approaches, as indicated by its lower sensitivity and FPR. The high FPR of the penalized likelihood based methods

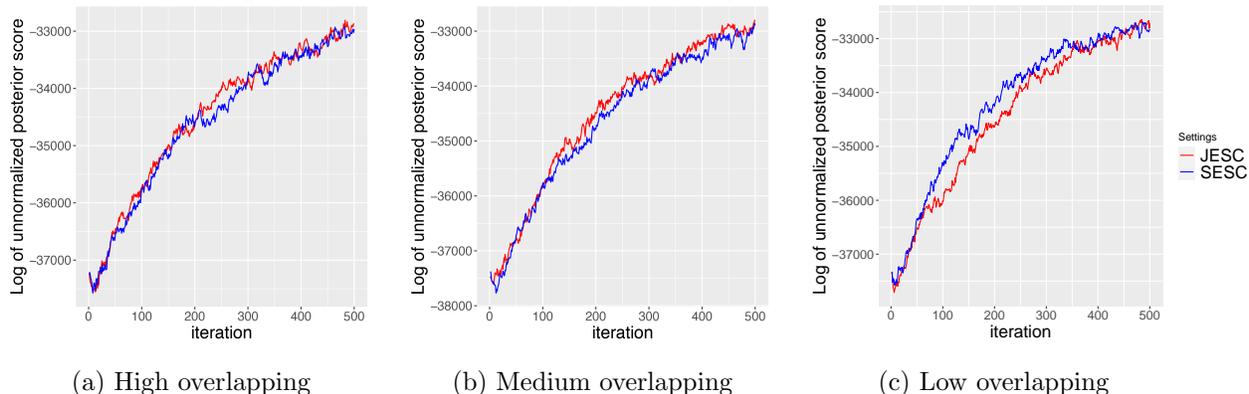| (a) High overlapping | (b) Medium overlapping | (c) Low overlapping |

Figure 1: The log of unnormalized posterior scores during the first 500 iterations under different scenarios.

in selecting differential edges is partly due to the fact that they select a larger number of false positive edges overall and may be because the regularization methods based on cross-validation tend to include many redundant variables resulting in a relatively larger number of errors compared with those for the Bayesian methods (Peterson et al.; 2020).

The proposed method achieves the highest MCC in identifying all edges across methods compared and yields a higher AUC compared with the separate inference, especially in the high and medium overlapping scenarios. As indicated in our theoretical results, the estimation performance based on the joint inference benefits the most when all graphs share the common support. Figure 1 shows the unnormalized posterior scores in log scale. Based on Figure 1, it seems that, in the high and medium overlapping settings, not only does JESC outperform SESC but also the posterior probabilities based on JESC increase faster than SESC during the beginning of the MCMC procedure.

## 5 Inferring Brain Functional Networks

In this section, we continue the illustration of JESC by applying the proposed method to an fMRI data set for simultaneously inferring multiple brain functional networks. Parkinson's disease (PD) is a major neurodegenerative disease influenced by both genetic and environmental factors (Halliday et al.; 2014). As the second most common neurodegenerative disorder, PD is characterized by the degeneration of dopamine-producing cells in the brain resulting in motor symptoms and nonmotor features (Mhyre et al.; 2012). Depression is the most common psychiatric symptom in patients with PD, and one of the earliest prodromal comorbidities that can have a significant impact on the quality of life (Chagas et al.; 2013). Nonmotor features including depression can appear in the earliest phase of the disease even before clinical motor impairment (Lix et al.; 2010; Shearer

15

et al.; 2012; Tibar et al.; 2018), but the efficacy of medications and psychotherapies for treating depression in PD (DPD) patients remains limited (Abós et al.; 2017). Hence, advances in timely detection and concerted management of DPD becomes urgent.

Up until now, the neural and pathophysiologic mechanisms of DPD remain unclear and are key research priorities for neurologists. A variety of neuroimaging technologies including fMRI, structure MRI, positron emission tomography and electroencephalography have been adopted to study PD. Among these, neuroimaging indicators have achieved considerable progress, and have provided new insights into PD. Resting-state fMRI exploits blood oxygen level-dependent signal to assess the correlation of the networks in different brain areas. An intra- and inter-network functional connectivity study in DPD demonstrated abnormal functional connection in left frontoparietal network, basal ganglia network, salience network and default-mode network (Wei et al.; 2017). To understand the underlying functional network changes for both DPD and non-depressed PD (NDPD) patients so that physicians could get an early-diagnosis in time for available treatment, we apply the proposed method to an fMRI data set (Wei et al.; 2017) for identifying regions of interest that are associated with the aberrant functional network and relevant to the onset of DPD and NDPD.

Twenty-one DPD patients, 49 NDPD patients and 50 matched healthy controls (HC) were recruited. Image data were acquired using a Siemens 3.0-Tesla signal scanner and functional imaging data were collected transversely by using a gradient-recalled echo-planar imaging (GRE-EPI) pulse sequence. We further perform image preprocessing procedure using Data Processing Assistant for Resting-State fMRI (http://rfmri.org/DPARSF) based on Statistical Parametric Mapping (SPM12, http://www.fil.ion.ucl.ac.uk/spm/) operated on the Matlab platform. Zang et al. (2004) proposed the method of Regional Homogeneity (ReHo) to analyze characteristics of regional brain activity and to reflect the temporal homogeneity of neural activity. In particular, we focus on the mReHo maps obtained by dividing the mean ReHo of the whole brain within each voxel in the ReHo map. We further segment the mReHo maps based on the Harvard-Oxford atlas (HOA) and extract all the mReHo signals corresponding to 15 subcortical regions of interest (ROI) (HOA number: 97-112) using the Resting-State fMRI Data Analysis Toolkit. Hence, adapted to our setting, $n_1 = 21$, $n_2 = 49$, $n_3 = 50$, $p = 15$, and the ordering is taken according to the HOA number.

We apply JESC along with other contenders to the resulting mReHo data set consisting of three groups for jointly estimating the functional connectivity networks. The parameter configuration are identical to those in the simulation study. Table 4 lists the number of edges selected by JESC and its competitors. The separate estimation methods (SESC and DAGL) resulted in graphs that share fewer edges in the Cholesky factors for the precision matrices of three groups. JGL resulted in most shared edges, followed by MPenPC and our method (JESC). Overall, JGL and MPenPC selected a lot more linked genes than other methods. JESC and SESC selected less unique edges among the ROIs for DPD than those for NDPD and HC. This might suggest that the patients with DPD lack

some important links among the subcortical regions. By visualizing the brain connectome as nodes and edges, Figure 2 shows the DAGs for three groups and all the shared edges estimated by JESC.

Table 5 lists six edges that are unique to the group of DPD identified by JESC. In particular, we discover discriminative connectivity changes between hippocampus and amygdala areas. These findings suggest disease-related alterations of functional connectivity as the basis for faulty information processing in DPD. Our findings are in good agreement with the aberrant functional features in subcortical regions that are related to the onset of DPD as shown in previous studies (Dan et al.; 2017; Lin et al.; 2020; Cao et al.; 2020).



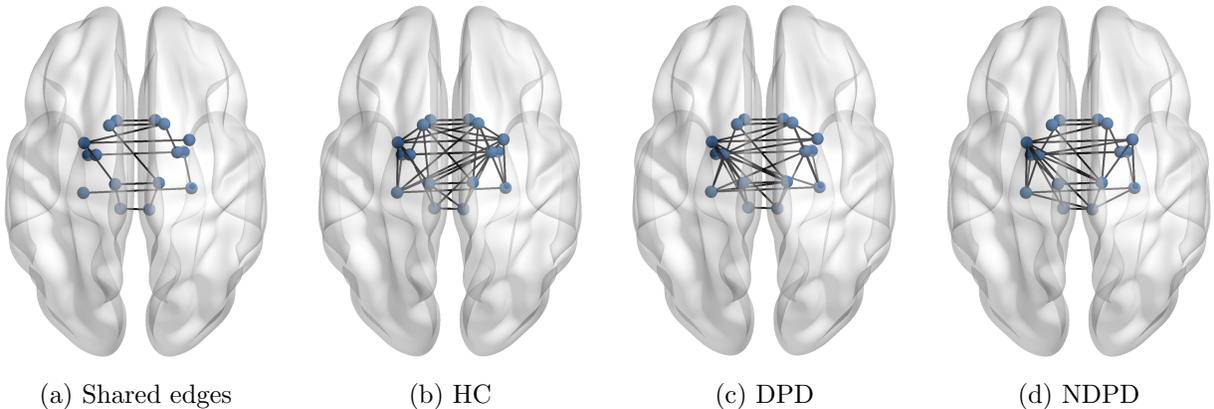|         (a) Shared edges          |          (b) HC          |          (c) DPD          |          (d) NDPD          |

Figure 2: Estimated brain function activity networks for HC, DPD, NDPD and the shared connections among three groups.

Table 4: Number of edges selected by the proposed method and its competitors. "DPD unique" counts the number of edges that only appear in the DPD group; "NDPD unique" counts the number of edges that only appear in the NDPD group; "HC unique" counts the number of edges that only appear in the HC group; and "Shared" counts the number of edges shared by all three groups.

| Method | DPD unique | NDPD unique | HC unique | Shared |
|--------|------------|-------------|-----------|--------|
| JESC   | 6          | 8           | 12        | 14     |
| SESC   | 8          | 9           | 11        | 10     |
| MPenPC | 10         | 5           | 8         | 19     |
| JGL    | 9          | 3           | 9         | 27     |
| DAGL   | 5          | 3           | 5         | 5      |

# 6   Discussion

In this paper, we proposed the JESC prior for Bayesian joint inference of multiple DAGs. In high-dimensional settings, the induced posterior attains the joint selection consistency under mild conditions. We also showed the advantage of the joint inference over separate inferences, in terms of

Table 5: Estimated edges that are unique to the group of DPD and the related brain regions indexed in the HOA template.

| ID | HOA number | Brain region A | HOA number | Brain region B |
|----|-----------|----------------|-----------|----------------|
| 1 | 105 | Left Pallidum | 99 | Left Thalamus |
| 2 | 105 | Left Pallidum | 102 | Right Caudate |
| 3 | 107 | Left Hippocampus | 100 | Right Thalamus |
| 4 | 107 | Left Hippocampus | 101 | Left Caudate |
| 5 | 108 | Right Hippocampus | 106 | Right Pallidum |
| 6 | 109 | Left Amygdala | 108 | Right Hippocampus |

requiring weaker beta-min conditions, when the DAGs share the common structure. The proposed joint inference outperforms other state-of-the-art methods in numerical studies based on simulated data sets. We also applied our method to an fMRI data set, where our results are consistent with previous neurological findings.

Throughout the paper, we focus on the MRF prior to encourage similar structures across all DAGs. The other choice of prior can be imposed that depends on the relationship between graphs. For example, if there is a natural ordering between $K$ classes so that it is expected that the DAGs were generated based on a Markov chain, one can use a prior,

$$
\begin{aligned}
f(S_{1j}, \ldots, S_{Kj}) &= f(S_{1j}) \prod_{k=2}^{K} \pi(S_{kj} \mid S_{k-1,j}) \\
&\propto \prod_{k=2}^{K} \exp\left\{ 2c_2 \sum_{l=1}^{j-1} I(S_{k,jl} = S_{k-1,jl} = 1) \right\}, \quad j = 2, \ldots, p
\end{aligned}
$$

for some constant $c_2 > 0$, which encourages similar patterns of sparsity for two consecutive graphs $S_{kj}$ and $S_{k-1,j}$. Theoretical properties of the joint inference based on various types of joint priors for $(S_{1j}, \ldots, S_{Kj})$, including the above Markov-type prior, may worth investigating as future work.

# 7 Proofs

**Proof of Theorem 3.1** Let $S_j = (S_{1j}, \ldots, S_{Kj})$ and $S_{0j} = (S_{01,j}, \ldots, S_{0K,j})$. It suffices to show that

$$\mathbb{E}_0\left\{\pi_\alpha\left(S_j \neq S_{0j} \mid \tilde{\mathbf{X}}_n\right)\right\} = o(p^{-1})$$

for any $j = 2, \ldots, p$, because

$$1 - \mathbb{E}_0\left\{\pi_\alpha\left(S_{A_1} = S_{A_{01}}, \ldots, S_{A_K} = S_{A_{0K}} \mid \tilde{\mathbf{X}}_n\right)\right\} \leq \sum_{j=2}^{p} \mathbb{E}_0\left\{\pi_\alpha\left(S_j \neq S_{0j} \mid \tilde{\mathbf{X}}_n\right)\right\}.$$

Note that $\{S_j \neq S_{0j}\}$ is equivalent to $\{S_{kj} \neq S_{0k,j}$ for at least one $k = 1, \ldots, K\}$. For given $1 \leq l \leq K$ and $1 \leq k_1 < \ldots < k_l \leq K$, define

$$N_{k_1,\ldots,k_l} := \left\{S_j : S_{kj} \neq S_{0k,j} \text{ if and only if } k \in \{k_1, \ldots, k_l\}\right\}.$$

Then, we have

$$\pi_\alpha\left(S_j \neq S_{0j} \mid \tilde{\mathbf{X}}_n\right)$$

$$= \sum_{k=1}^{K} \sum_{S_j \in N_k} \pi_\alpha\left(S_j \mid \tilde{\mathbf{X}}_n\right) + \sum_{k_1 < k_2} \sum_{S_j \in N_{k_1,k_2}} \pi_\alpha\left(S_j \mid \tilde{\mathbf{X}}_n\right) + \sum_{k_1 < k_2 < k_3} \sum_{S_j \in N_{k_1,k_2,k_3}} \pi_\alpha\left(S_j \mid \tilde{\mathbf{X}}_n\right)$$

$$+ \cdots + \sum_{S_j \in N_{1,\ldots,K}} \pi_\alpha\left(S_j \mid \tilde{\mathbf{X}}_n\right). \tag{8}$$

The first term in (8) can be divided into two parts:

$$\sum_{k=1}^{K} \mathbb{E}_0\left\{\pi_\alpha(S_j \in N_k \mid \tilde{\mathbf{X}}_n)\right\}$$

$$= \sum_{k=1}^{K} \left[\mathbb{E}_0\left\{\pi_\alpha(S_j \in N_k, S_{kj} \supseteq S_{0k,j} \mid \tilde{\mathbf{X}}_n)\right\} + \mathbb{E}_0\left\{\pi_\alpha(S_j \in N_k, S_{kj} \not\supseteq S_{0k,j} \mid \tilde{\mathbf{X}}_n)\right\}\right].$$

Let $\pi^I(S_{kj} \mid \mathbf{X}_{n_k}) \propto f_\alpha(\mathbf{X}_{n_k} \mid S_{kj})\pi(S_{kj})$ be the posterior for $S_{kj}$ based on the separate inference for each DAG. Note that if $S_j \in N_k$, then we have

$$\frac{\pi_\alpha\left(S_j \mid \tilde{\mathbf{X}}_n\right)}{\pi_\alpha\left(S_{0j} \mid \tilde{\mathbf{X}}_n\right)} = \frac{\pi_\alpha^I(S_{kj} \mid \mathbf{X}_{n_k})}{\pi_\alpha^I(S_{0k,j} \mid \mathbf{X}_{n_k})} \frac{f(S_{1j}, \ldots, S_{Kj})}{f(S_{01,j}, \ldots, S_{0K,j})}$$

and

$$\frac{f(S_{1j}, \ldots, S_{Kj})}{f(S_{01,j}, \ldots, S_{0K,j})} = \exp\left[c_{2j} \sum_{k' \neq k} \left\{|S_{kj} \cap S_{0k',j}| - |S_{0k,j} \cap S_{0k',j}|\right\}\right]$$

$$\leq \exp\left[c_{2j} \sum_{k' \neq k} \left\{|S_{kj} \cap S_{0k',j}|\right\}\right]$$

$$\leq \exp\left\{c_{2j}(K-1)(j-1)\right\} \leq \exp\left\{c_{2j}(j-1)K\right\}.$$

19

Then, we have

$$\sum_{k=1}^{K} \Big[ \mathbb{E}_0 \big\{ \pi_\alpha(S_j \in N_k, S_{kj} \supsetneq S_{0k,j} \mid \tilde{\mathbf{X}}_n) \big\}$$

$$\leq \sum_{k=1}^{K} \sum_{S_j \in N_k, S_{kj} \supsetneq S_{0k,j}} \mathbb{E}_0 \Big\{ \frac{\pi_\alpha(S_j \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{0j} \mid \tilde{\mathbf{X}}_n)} \Big\}$$

$$= \sum_{k=1}^{K} \sum_{S_j \in N_k, S_{kj} \supsetneq S_{0k,j}} \mathbb{E}_0 \Big\{ \prod_{k'=1}^{K} \frac{\pi_\alpha^I(S_{k'j} \mid \mathbf{X}_{n_{k'}})}{\pi_\alpha^I(S_{0k',j} \mid \mathbf{X}_{n_{k'}})} \Big\} \frac{f(S_{01,j}, \ldots, S_{0k-1,j}, S_{kj}, S_{0k+1,j}, \ldots, S_{0K,j})}{f(S_{01,j}, \ldots, S_{0K,j})}$$

$$\leq \sum_{k=1}^{K} \sum_{S_{kj} \supsetneq S_{0k,j}} \mathbb{E}_0 \Big\{ \frac{\pi_\alpha^I(S_{kj} \mid \mathbf{X}_{n_k})}{\pi_\alpha^I(S_{0k,j} \mid \mathbf{X}_{n_k})} \Big\} \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\lesssim K p^{-c_1} R_j \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\leq K p^{-(c_1-1)} \exp\big\{ c_{2j}(j-1)K \big\} \leq K p^{-\{(C_{\mathrm{bm}}-c_1-1)\wedge(c_1-1)\}} \exp\big\{ c_{2j}(j-1)K \big\}$$

where the last inequality follows from the proof of Lemma 6.1 in Lee et al. (2019). Let $N_{S_{kj},\alpha,\chi^2}$ be the set defined in the proof of Theorem 3.1 in Lee et al. (2019). Then,

$$\sum_{k=1}^{K} \mathbb{E}_0 \big\{ \pi_\alpha(S_j \in N_k, S_{kj} \not\supseteq S_{0k,j} \mid \tilde{\mathbf{X}}_n) \big\} \Big]$$

$$\leq \sum_{k=1}^{K} \sum_{S_{kj} \not\supseteq S_{0k,j}} \mathbb{P}_0 \big( \mathbf{X}_{n_k} \in N_{S_{kj},\alpha,\chi^2} \big)$$

$$+ \sum_{k=1}^{K} \sum_{S_{kj} \not\supseteq S_{0k,j}} \mathbb{E}_0 \Big\{ \frac{\pi_\alpha^I(S_{kj} \mid \mathbf{X}_{n_k})}{\pi_\alpha^I(S_{0k,j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N_{S_{kj},\alpha,\chi^2}^c) \Big\} \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\lesssim \sum_{k=1}^{K} \sum_{S_{kj} \not\supseteq S_{0k,j}} \exp\Big\{ -\frac{(\epsilon')^2 \epsilon_0^2}{64(1+2\epsilon_0)^2} n_k \Big\}$$

$$+ \sum_{k=1}^{K} \sum_{S_{kj} \not\supseteq S_{0k,j}} \mathbb{E}_0 \Big\{ \frac{\pi_\alpha^I(S_{kj} \mid \mathbf{X}_{n_k})}{\pi_\alpha^I(S_{0k,j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N_{S_{kj},\alpha,\chi^2}^c) \Big\} \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\lesssim K \exp\Big\{ -\frac{(\epsilon')^2 \epsilon_0^2}{128(1+2\epsilon_0)^2} \min_k n_k \Big\} + K \big( p^{-C_{\mathrm{bm}}+1} R_j + p^{-C_{\mathrm{bm}}+c_1+1} \big) \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\lesssim K p^{-(C_{\mathrm{bm}}-c_1-1)} \exp\big\{ c_{2j}(j-1)K \big\}$$

$$\leq K p^{-\{(C_{\mathrm{bm}}-c_1-1)\wedge(c_1-1)\}} \exp\big\{ c_{2j}(j-1)K \big\}$$

where the second and third inequalities follow from Lemma 6.2 and the proof of Theorem 3.1 in Lee et al. (2019). The last inequality holds by Condition (P) because we assume $s_0 \geq C_{\mathrm{bm}} - c_1 - 1$.

Now consider the second term in (8). Note that if $S_j \in N_{k_1,k_2}$, then we have

$$\frac{\pi_\alpha(S_j \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{0j} \mid \tilde{\mathbf{X}}_n)} = \frac{\pi_\alpha^I(S_{k_1j} \mid \mathbf{X}_{n_{k_1}})}{\pi_\alpha^I(S_{0k_1,j} \mid \mathbf{X}_{n_{k_1}})} \frac{\pi_\alpha^I(S_{k_2j} \mid \mathbf{X}_{n_{k_2}})}{\pi_\alpha^I(S_{0k_2,j} \mid \mathbf{X}_{n_{k_2}})} \frac{f(S_{1j}, \ldots, S_{Kj})}{f(S_{01,j}, \ldots, S_{0K,j})}$$

and

$$\frac{f(S_{1j}, \ldots, S_{Kj})}{f(S_{01,j}, \ldots, S_{0K,j})} \leq \exp\left[c_{2j}\left\{|S_{k_1j} \cap S_{k_2j}| + \sum_{k' \notin \{k_1, k_2\}} |S_{klj} \cap S_{0k',j}|\right\}\right]$$
$$\leq \exp\left[c_{2j}\{j - 1 + 2(K-2)(j-1)\}\right]$$
$$\leq \exp\{c_{2j}(j-1)2K\}.$$

If $S_j \in N_{k_1,k_2}$, then one of the followings holds: (1) $S_{k_1j} \supsetneq S_{0k_1,j}$ and $S_{k_2j} \supsetneq S_{0k_2,j}$, (2) $S_{k_1j} \supsetneq S_{0k_1,j}$ and $S_{k_2j} \not\supseteq S_{0k_2,j}$, (3) $S_{k_1j} \not\supseteq S_{0k_1,j}$ and $S_{k_2j} \supsetneq S_{0k_2,j}$ or (4) $S_{k_1j} \not\supseteq S_{0k_1,j}$ and $S_{k_2j} \not\supseteq S_{0k_2,j}$. For example, by the similar arguments used in the previous paragraph,

$$\sum_{k_1 < k_2} \sum_{\substack{S_j \in N_{k_1,k_2}, \\ S_{k_1j} \supsetneq S_{0k_1,j}, S_{k_2j} \not\supseteq S_{0k_2,j}}} \mathbb{E}_0\{\pi_\alpha(S_j \mid \tilde{\mathbf{X}}_n)\}$$

$$\lesssim \sum_{k_1 < k_2} \sum_{S_{k_2j} \not\supseteq S_{0k_2,j}} \mathbb{P}_0\left(\mathbf{X}_{n_{k_2}} \in N_{S_{k_2j},\alpha,\chi^2}\right)$$

$$+ \sum_{k_1 < k_2} \sum_{S_{k_1j} \supsetneq S_{0k_1,j}} \mathbb{E}_0\left\{\frac{\pi_\alpha^I(S_{k_1j} \mid \mathbf{X}_{n_{k_1}})}{\pi_\alpha^I(S_{0k_1,j} \mid \mathbf{X}_{n_{k_1}})}\right\} \sum_{S_{k_2j} \not\supseteq S_{0k_2,j}} \mathbb{E}_0\left\{\frac{\pi_\alpha^I(S_{k_2j} \mid \mathbf{X}_{n_{k_2}})}{\pi_\alpha^I(S_{0k_2,j} \mid \mathbf{X}_{n_{k_2}})} I(\mathbf{X}_{n_{k_2}} \in N_{S_{k_2j},\alpha,\chi^2}^c)\right\}$$

$$\times \exp\{c_{2j}(j-1)2K\}$$

$$\lesssim \sum_{k_1 < k_2} p^{-2\{(C_{\text{bm}}-c_1-1)\wedge(c_1-1)\}} \exp\{c_{2j}(j-1)2K\}$$

$$\leq K^2 p^{-2\{(C_{\text{bm}}-c_1-1)\wedge(c_1-1)\}} \exp\{c_{2j}(j-1)2K\}.$$

Thus, by applying the similar arguments to the above four cases, the expectation of the second term in (8) is

$$\sum_{k_1 < k_2} \sum_{S_j \in N_{k_1,k_2}} \mathbb{E}_0\{\pi_\alpha(S_j \mid \tilde{\mathbf{X}}_n)\}$$

$$\leq \sum_{k_1 < k_2} \sum_{S_j \in N_{k_1,k_2}} \mathbb{E}_0\left\{\frac{\pi_\alpha^I(S_j \mid \mathbf{X}_n)}{\pi_\alpha^I(S_{0,j} \mid \mathbf{X}_n)}\right\} \frac{f(S_j)}{f(S_{0j})}$$

$$\leq \sum_{k_1 < k_2} \sum_{S_{k_1j} \neq S_{0k_1,j}} \sum_{S_{k_2j} \neq S_{0k_2,j}} \mathbb{E}_0\left\{\frac{\pi_\alpha^I(S_{k_1j} \mid \mathbf{X}_{n_{k_1}})}{\pi_\alpha^I(S_{0k_1,j} \mid \mathbf{X}_{n_{k_1}})} \frac{\pi_\alpha^I(S_{k_2j} \mid \mathbf{X}_{n_{k_2}})}{\pi_\alpha^I(S_{0k_2,j} \mid \mathbf{X}_{n_{k_2}})}\right\} \exp\{c_{2j}(j-1)2K\}$$

$$\lesssim K^2 p^{-2\{(C_{\text{bm}}-c_1-1)\wedge(c_1-1)\}} \exp\{c_{2j}(j-1)2K\}.$$

Note that if $S_j \in N_{k_1,\ldots,k_l}$, then we have

$$\frac{f(S_{1j}, \ldots, S_{Kj})}{f(S_{01,j}, \ldots, S_{0K,j})} \leq \exp\left[c_{2j}\left\{\frac{l(l-1)}{2}(j-1) + l(K-l)(j-1)\right\}\right]$$
$$\leq \exp\{c_{2j}(j-1)lK\}.$$

Therefore, by repeatedly applying the similar arguments, we have

$$
\begin{aligned}
\mathbb{E}_0\Big\{\pi_\alpha\Big(S_j \neq S_{0j} \mid \tilde{\mathbf{X}}_n\Big)\Big\} &\lesssim \sum_{k=1}^{K}\Big[\frac{K\exp\{c_{2j}(j-1)K\}}{p^{\{(C_{\mathrm{bm}}-c_1-1)\wedge(c_1-1)\}}}\Big]^k \\
&\leq \sum_{k=1}^{K}\Big[\frac{Ke^K}{p^{\{(C_{\mathrm{bm}}-c_1-1)\wedge(c_1-1)\}}}\Big]^k \\
&\lesssim \frac{Ke^K}{p^{\{(C_{\mathrm{bm}}-c_1-1)\wedge(c_1-1)\}}} = o(p^{-1}),
\end{aligned}
$$

because we assume $C_{\mathrm{bm}} > c_1 + 2$, $c_1 > 2$, $c_{2j} \leq 1/(j-1)$ and $K = o(\log p)$.

**Proof of Theorem 3.2** We will only show that (6) holds for $k = 1$ when $\cup_{k'=2}^{K}S_{0k',j} \subseteq S_{01,j}$, but one can easily check the other cases using similar arguments. Let $S_j = (S_{1j}, \ldots, S_{Kj})$ and $S_{0j} = (S_{01,j}, \ldots, S_{0K,j})$. Note that

$$
\frac{\pi_\alpha(S_j \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{0j} \mid \tilde{\mathbf{X}}_n)} = \frac{\pi_\alpha^I(S_j \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha^I(S_{0j} \mid \tilde{\mathbf{X}}_n)}\frac{f(S_j)}{f(S_{0j})}.
$$

Because we assume $\cup_{k'=2}^{K}S_{0k',j} \subseteq S_{01,j}$, it holds that $f(S_{1,j}, S_{02,j}, \ldots, S_{0K,j}) \leq f(S_{01,j}, S_{02,j}, \ldots, S_{0K,j})$ for any $S_{1j} \neq S_{01,j}$. Thus,

$$
\begin{aligned}
\frac{\pi_\alpha(S_{1j}, S_{02,j}, \ldots, S_{0K,j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{01,j}, S_{02,j}, \ldots, S_{0K,j} \mid \tilde{\mathbf{X}}_n)} &= \frac{\pi_\alpha^I(S_{1j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha^I(S_{01,j} \mid \tilde{\mathbf{X}}_n)}\frac{f(S_{1j}, S_{02,j}, \ldots, S_{0K,j})}{f(S_{01,j}, S_{02,j}, \ldots, S_{0K,j})} \\
&\leq \frac{\pi_\alpha^I(S_{1j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha^I(S_{01,j} \mid \tilde{\mathbf{X}}_n)}
\end{aligned}
$$

for any $S_{1j} \neq S_{01,j}$. Then, we have

$$
\begin{aligned}
\frac{1 - \pi_\alpha(S_{01,j} \mid S_{02,j}, \ldots, S_{0K,j}, \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{01,j} \mid S_{02,j}, \ldots, S_{0K,j}\tilde{\mathbf{X}}_n)} &= \sum_{S_{1j} \neq S_{01,j}} \frac{\pi_\alpha(S_{1j} \mid S_{02,j}, \ldots, S_{0K,j}, \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{01,j} \mid S_{02,j}, \ldots, S_{0K,j}\tilde{\mathbf{X}}_n)} \\
&= \sum_{S_{1j} \neq S_{01,j}} \frac{\pi_\alpha(S_{1j}, S_{02,j}, \ldots, S_{0K,j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha(S_{01,j}, S_{02,j}, \ldots, S_{0K,j} \mid \tilde{\mathbf{X}}_n)} \\
&\leq \sum_{S_{1j} \neq S_{01,j}} \frac{\pi_\alpha^I(S_{1j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha^I(S_{01,j} \mid \tilde{\mathbf{X}}_n)} = \frac{1 - \pi_\alpha^I(S_{01,j} \mid \tilde{\mathbf{X}}_n)}{\pi_\alpha^I(S_{01,j} \mid \tilde{\mathbf{X}}_n)},
\end{aligned}
$$

which implies

$$
\pi_\alpha\Big(S_{01,j} \mid S_{01,j}, \ldots, S_{02,j}, \ldots, S_{0K,j}, \tilde{\mathbf{X}}_n\Big) \geq \pi_\alpha^I(S_{01,j} \mid \mathbf{X}_{n_k}).
$$

**Proof of Theorem 3.3** In this proof, let $S_j = (S_{1j}, \ldots, S_{j-1j})$ and $S_{0j} = (S_{0,1j}, \ldots, S_{0,j-1j})$ be the (common) support of the $j$th row of the lower triangular part of $S_A$ and $S_0$, respectively. Let

$$
\tilde{\pi}_\alpha(S_A \mid \tilde{\mathbf{X}}_n) \propto \pi_\alpha(S_{A_1} = \cdots = S_{A_K} = S_A \mid \tilde{\mathbf{X}}_n)
$$

be the joint posterior for $(S_{A_1}, \ldots, S_{A_K})$ restricted to common supports. Then,

$$
\begin{aligned}
\tilde{\pi}_\alpha(S_A \neq S_0 \mid \tilde{\mathbf{X}}_n) &\leq \sum_{j=2}^{p} \tilde{\pi}_\alpha(S_j \neq S_{0j} \mid \tilde{\mathbf{X}}_n) \\
&= \sum_{j=2}^{p} \tilde{\pi}_\alpha(S_j \supsetneq S_{0j} \mid \tilde{\mathbf{X}}_n) + \sum_{j=2}^{p} \tilde{\pi}_\alpha(S_j \not\supseteq S_{0j} \mid \tilde{\mathbf{X}}_n).
\end{aligned} \tag{9}
$$

The expectation of the first part of (9) is bounded above by

$$
\begin{aligned}
\sum_{j=2}^{p} \mathbb{E}_0 \Big\{ \tilde{\pi}_\alpha(S_j \supsetneq S_{0j} \mid \tilde{\mathbf{X}}_n) \Big\} &\leq \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \mathbb{E}_0 \Big\{ \frac{\tilde{\pi}_\alpha(S_j \mid \tilde{\mathbf{X}}_n)}{\tilde{\pi}_\alpha(S_{0j} \mid \tilde{\mathbf{X}}_n)} \Big\} \\
&= \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \mathbb{E}_0 \Big\{ \prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^I(S_{kj} = S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^I(S_{kj} = S_{0j} \mid \mathbf{X}_{n_k})} \Big\} \frac{f(S_j, \ldots, S_j)}{f(S_{0j}, \ldots, S_{0j})} \\
&\lesssim \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \Big\}^K c_{\alpha,\gamma}^{K(|S_j| - |S_{0j}|)} \exp \big\{ c_{2j} K(K-1)(|S_j| - |S_{0j}|) \big\} \\
&\lesssim \sum_{j=2}^{p} c_{\alpha,\gamma}^K p^{-c_1 K} R_j \exp \big\{ c_{2j} K(K-1)(j-1) \big\} \\
&\lesssim \exp \Big\{ -c_1 K \log p + \log p + K \log c_{\alpha,\gamma} + K(K-1) \Big\} = o(1),
\end{aligned}
$$

where $c_{\alpha,\gamma} = (1 + \alpha/\gamma)^{-1/2} \{2/(1-\alpha)\}^{1/2}$, by the proof of Lemma 6.1 in Lee et al. (2019), $c_1 > 1$, $c_{2j} \leq 1/(j-1)$ and $K = o(\log p)$.

On the other hand, the expectation of the second part of (9) is

$$
\begin{aligned}
&\sum_{j=2}^{p} \mathbb{E}_0 \Big\{ \tilde{\pi}_\alpha(S_j \not\supseteq S_{0j} \mid \tilde{\mathbf{X}}_n) \Big\} \\
&\leq \sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \sum_{k=1}^{K} \mathbb{P}_0\big(\mathbf{X}_{n_k} \in N_{S_j, \alpha, \chi^2}\big) \\
&\quad + \sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \mathbb{E}_0 \Big\{ \tilde{\pi}_\alpha(S_j \not\supseteq S_{0j} \mid \tilde{\mathbf{X}}_n) I(\mathbf{X}_{n_k} \in N_{S_j, \alpha, \chi^2}^c, \forall k) \Big\} \\
&\lesssim pK \exp \Big\{ -\frac{(\epsilon')^2 \epsilon_0^2}{128(1 + 2\epsilon_0)^2} \min_k n_k \Big\} \\
&\quad + \sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \mathbb{E}_0 \Big\{ \prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^I(S_{kj} = S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^I(S_{kj} = S_{0j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N_{S_j, \alpha, \chi^2}^c) \Big\} \frac{f(S_j, \ldots, S_j)}{f(S_{0j}, \ldots, S_{0j})}.
\end{aligned}
$$

(10)

(11)

Note that (10) is of order $o(1)$ and

$$
\frac{f(S_j, \ldots, S_j)}{f(S_{0j}, \ldots, S_{0j})} \leq \exp \big\{ c_{2j} K(K-1) ||S_j| - |S_{0j}|| \big\} \leq \exp \big\{ K(K-1) \big\}.
$$

23

By the proof of Theorem 3.1 in Lee et al. (2019),

$$
\mathbb{E}_0\Big\{ \prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^I(S_{kj}=S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^I(S_{kj}=S_{0j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N^c_{S_j,\alpha,\chi^2}) \Big\}
$$

$$
\leq \prod_{k=1}^{K} \frac{\pi(S_j)}{\pi(S_{0j})} \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \exp\Big\{ -\frac{\alpha(1-\alpha)}{4}\frac{\epsilon_0^2(1-2\epsilon_0)^2}{4} n_k \|a_{0k,S_{0j}\cap S_j^c}\|_2^2 \Big\}
$$

$$
+\Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \Big\}^K \sum_{k=1}^{K} \mathbb{P}_0(\mathbf{X}_{n_k} \in N_{j,S_{kj}})
$$

$$
\leq \Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \Big\}^K \exp\Big\{ -(|S_{0j}|-|S_j\cap S_{0j}|)C_{\mathrm{bm}}K\log p \Big\}
$$

$$
+\Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \Big\}^K \sum_{k=1}^{K} 4\exp\big(-n_k\epsilon_0^2/2\big)
$$

$$
\lesssim \Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \Big\}^K \Big\{ \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \Big\}^K \exp\Big\{ -(|S_{0j}|-|S_j\cap S_{0j}|)C_{\mathrm{bm}}K\log p \Big\}
$$

where $N_{j,S_{kj}}$ is the set defined in the proof of Theorem 3.1 in Lee et al. (2019), $\nu_1 = (1+\alpha/\gamma)^{1/2}$ and $\nu_2 = \{1-(\alpha+\nu_0/n)/(1-4\sqrt{\epsilon'}-5\epsilon')\}^{-1/2}$. Note that the last inequality holds due to $K\log p = o(\min_k n_k)$ and

$$
\exp\Big\{ -\frac{\alpha(1-\alpha)}{4}\frac{\epsilon_0^2(1-2\epsilon_0)^2}{4} \sum_{k=1}^{K} n_k \|a_{0k,S_{0j}\cap S_j^c}\|_2^2 \Big\}
$$

$$
\leq \exp\Big\{ -\frac{\alpha(1-\alpha)}{4}\frac{\epsilon_0^2(1-2\epsilon_0)^2}{4}(|S_{0j}|-|S_j\cap S_{0j}|)\min_{(j,l):a_{01,jl}\neq 0}\sum_{k=1}^{K} n_k a_{0k,jl}^2 \Big\}
$$

$$
\leq \exp\Big\{ -(|S_{0j}|-|S_j\cap S_{0j}|)C_{\mathrm{bm}}K\log p \Big\}
$$

for some constant $C>0$ by condition (B3). Thus, (11) is bounded above by

$$
\sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \Big\{ \frac{\pi(S_j)}{\pi(S_{0j})} \Big\}^K \Big\{ \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j}\cap S_j|} \Big\}^K
$$

$$
\times \exp\Big\{ -(|S_{0j}|-|S_j\cap S_{0j}|)C_{\mathrm{bm}}K\log p + K(K-1) \Big\}
$$

$$
\lesssim \exp\Big\{ -(C_{\mathrm{bm}}-c_1-2)K\log p + K(K-1) \Big\} = o(1),
$$

because $C_{\mathrm{bm}} > c_1 + 2$ and $K = o(\log p)$. This completes the proof.

**Proof of Theorem 3.4** Similarly to the proof of Theorem 3.3, we have

$$
\sum_{j=2}^{p} \mathbb{E}_0\Big\{\tilde{\pi}_\alpha^*(S_j \supsetneq S_{0j} \mid \tilde{\mathbf{X}}_n)\Big\} \leq \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \mathbb{E}_0\Big\{\prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_{0j} \mid \mathbf{X}_{n_k})}\Big\} \frac{\tilde{f}(S_j,\ldots,S_j)}{\tilde{f}(S_{0j},\ldots,S_{0j})}
$$

$$
\lesssim \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \Big\{\frac{\tilde{\pi}(S_j)}{\tilde{\pi}(S_{0j})}\Big\}^K c_{\alpha,\gamma}^{K(|S_j|-|S_{0j}|)} \exp(K-1)
$$

$$
= \sum_{j=2}^{p} \sum_{S_j \supsetneq S_{0j}} \Big\{\frac{\pi(S_j)}{\pi(S_{0j})}\Big\} c_{\alpha,\gamma}^{K(|S_j|-|S_{0j}|)} \exp(K-1)
$$

$$
\lesssim \sum_{j=2}^{p} c_{\alpha,\gamma}^{K} p^{-c_1} R_j \exp(K-1)
$$

$$
\lesssim \exp\Big\{-(c_1-1)\log p + K \log c_{\alpha,\gamma} + K\Big\} = o(1),
$$

where $c_{\alpha,\gamma} = (1 + \alpha/\gamma)^{-1/2}\{2/(1-\alpha)\}^{1/2}$ and $\pi_\alpha^{I,*}(S_{kj} \mid \mathbf{X}_{n_k}) \propto f(\mathbf{X}_{n_k} \mid S_{kj})\tilde{\pi}(S_{kj})$, because $c_1 > 1$, $c_{2j} \leq 1/(j-1)$ and $K = o(\log p)$. The second and third inequalities hold by the proof of Lemma 6.1 in Lee et al. (2019) and

$$
\frac{\tilde{f}(S_j,\ldots,S_j)}{\tilde{f}(S_{0j},\ldots,S_{0j})} \leq \exp\big\{c_{2j}(K-1)\big||S_j|-|S_{0j}|\big|\big\} \leq \exp(K-1).
$$

Furthermore,

$$
\sum_{j=2}^{p} \mathbb{E}_0\Big\{\tilde{\pi}_\alpha^*(S_j \not\supseteq S_{0j} \mid \tilde{\mathbf{X}}_n)\Big\}
$$

$$
\lesssim pK \exp\Big\{-\frac{(\epsilon')^2 \epsilon_0^2}{128(1+2\epsilon_0)^2} \min_k n_k\Big\}
$$

$$
+ \sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \mathbb{E}_0\Big\{\prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_{0j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N_{S_j,\alpha,\chi^2}^c)\Big\} \frac{\tilde{f}(S_j,\ldots,S_j)}{\tilde{f}(S_{0j},\ldots,S_{0j})}, \quad (12)
$$

where

$$
\mathbb{E}_0\Big\{\prod_{k=1}^{K} \frac{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_j \mid \mathbf{X}_{n_k})}{\tilde{\pi}_\alpha^{I,*}(S_{kj} = S_{0j} \mid \mathbf{X}_{n_k})} I(\mathbf{X}_{n_k} \in N_{S_j,\alpha,\chi^2}^c)\Big\}
$$

$$
\lesssim \frac{\pi(S_j)}{\pi(S_{0j})}\Big\{\nu_1^{|S_{0j}|-|S_j|}\nu_2^{|S_j|-|S_{0j}\cap S_j|}\Big\}^K \exp\Big\{-(|S_{0j}|-|S_j\cap S_{0j}|)C_{\mathrm{bm}}\log p\Big\},
$$

by the similar arguments used in the proof of Theorem 3.3 and condition (C3). Note that the last inequality holds due to $\log p = o(\min_k n_k)$.

25

Therefore, (12) is bounded above by

$$\sum_{j=2}^{p} \sum_{S_j \not\supseteq S_{0j}} \frac{\pi(S_j)}{\pi(S_{0j})} \Big\{ \nu_1^{|S_{0j}|-|S_j|} \nu_2^{|S_j|-|S_{0j} \cap S_j|} \Big\}^K$$

$$\times \exp\Big\{ -(|S_{0j}| - |S_j \cap S_{0j}|)C_{\mathrm{bm}} \log p + (K-1) \Big\}$$

$$\lesssim \exp\Big\{ -(C_{\mathrm{bm}} - c_1 - 2)\log p + 2K \Big\} = o(1),$$

because $C_{\mathrm{bm}} > c_1 + 2$ and $K = o(\log p)$. This completes the proof.

# References

Abós, A., Baggio, H. C., Segura, B., García-Díaz, A. I., Compta, Y., Martí, M. J., Valldeoriola, F. and Junqué, C. (2017). Discriminating cognitive status in parkinson's disease through functional connectomics and machine learning, *Scientific Reports* **7**(1): 45347.

Banerjee, S. and Ghosal, S. (2015). Bayesian structure learning in graphical models, *Journal of Multivariate Analysis* **136**: 147–162.

Ben-David, E., Li, T., Massam, H. and Rajaratnam, B. (2015). High dimensional bayesian inference for gaussian directed acyclic graph models, *arXiv:1109.4371v5* .

Bickel, P. J. and Levina, E. (2008). Regularized estimation of large covariance matrices, *The Annals of Statistics* **36**(1): 199–227.

Cai, T. T., Li, H., Liu, W. and Xie, J. (2016). Joint estimation of multiple high-dimensional precision matrices, *Statistica Sinica* **26**(2): 445–464.

Cai, T. T., Zhang, C.-H. and Zhou, H. H. (2010). Optimal rates of convergence for covariance matrix estimation, *The Annals of Statistics* **38**(4): 2118–2144.

Cai, T. T. and Zhou, H. H. (2012b). Optimal rates of convergence for sparse covariance matrix estimation, *The Annals of Statistics* **40**(5): 2389–2420.

Cao, X., Khare, K. and Ghosh, M. (2019). Posterior graph selection and estimation consistency for high-dimensional bayesian dag models, *The Annals of Statistics* **47**(1): 319–348.

Cao, X., Wang, X., Xue, C., Zhang, S., Huang, Q. and Liu, W. (2020). A radiomics approach to predicting parkinson's disease by incorporating whole-brain functional activity and gray matter structure, *Frontiers in Neuroscience* **14**: 751.

Chagas, M. H. N., Linares, I. M., Garcia, G. J., Hallak, J. E., Tumas, V. and Crippa, J. A. S. (2013). Neuroimaging of depression in parkinson's disease: a review, *International Psychogeriatrics* **25**(12): 1953–1961.

Chen, J. and Chen, Z. (2008). Extended Bayesian information criteria for model selection with large model spaces, *Biometrika* **95**(3): 759–771.

Dan, R., Růžička, F., Bezdicek, O., Růžička, E., Roth, J., Vymazal, J., Goelman, G. and Jech, R. (2017). Separate neural representations of depression, anxiety and apathy in parkinson's disease, *Scientific Reports* **7**(1): 12164.

Danaher, P., Wang, P. and Witten, D. M. (2014). The joint graphical lasso for inverse covariance estimation across multiple classes, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**(2): 373–397.

Gan, L., Yang, X., Narisetty, N. and Liang, F. (2019). Bayesian joint estimation of multiple graphical models, *Advances in Neural Information Processing Systems*, pp. 9802–9812.

Halliday, G. M., Leverenz, J. B., Schneider, J. S. and Adler, C. H. (2014). The neurobiological basis of cognitive impairment in Parkinson's disease, *Movement Disorders* **29**(5): 634–650.

Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems, *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* **186**(1007): 453–461.

Khare, K., Oh, S.-Y., Rahman, S. and Rajaratnam, B. (2019). A scalable sparse cholesky based approach for learning high-dimensional covariance matrices in ordered data, *Machine Learning* **108**(12): 2061–2086.

Lee, K. and Lee, J. (2017). Estimating large precision matrices via modified cholesky decomposition, *Statistica Sinica* (accepted).

Lee, K., Lee, J. and Lin, L. (2019). Minimax posterior convergence rates and model selection consistency in high-dimensional dag models based on sparse cholesky factors, *The Annals of Statistics* **47**(6): 3413–3437.

Lin, H., Cai, X., Zhang, D., Liu, J., Na, P. and Li, W. (2020). Functional connectivity markers of depression in advanced parkinson's disease, *NeuroImage: Clinical* **25**: 102130.

Liu, J., Sun, W. and Liu, Y. (2019). Joint skeleton estimation of multiple directed acyclic graphs for heterogeneous population, *Biometrics* **75**(1): 36–47.

Lix, L. M., Hobson, D. E., Azimaee, M., Leslie, W. D., Burchill, C. and Hobson, S. (2010). Socioeconomic variations in the prevalence and incidence of parkinson's disease: a population-based analysis, *Journal of Epidemiology & Community Health* **64**(4): 335–340.

Martin, R., Mess, R. and Walker, S. G. (2017). Empirical bayes posterior concentration in sparse high-dimensional linear models, *Bernoulli* **23**(3): 1822–1847.

Martin, R. and Walker, S. G. (2014). Asymptotically minimax empirical bayes estimation of a sparse normal mean vector, *Electronic Journal of Statistics* **8**(2): 2188–2206.

Mhyre, T. R., Boyd, J. T., Hamill, R. W. and Maguire-Zeiss, K. A. (2012). *Parkinson's Disease*, pp. 389–455.

Peterson, C. B., Osborne, N., Stingo, F. C., Bourgeat, P., Doecke, J. D. and Vannucci, M. (2020). Bayesian modeling of multiple structural connectivity networks during the progression of alzheimer's disease, *Biometrics, to appear* .

Peterson, C., Stingo, F. C. and Vannucci, M. (2015). Bayesian inference of multiple gaussian graphical models, *Journal of the American Statistical Association* **110**(509): 159–174.

Ren, Z., Sun, T., Zhang, C.-H. and Zhou, H. H. (2015). Asymptotic normality and optimalities in estimation of large gaussian graphical models, *The Annals of Statistics* **43**(3): 991–1026.

Shearer, J., Green, C., Counsell, C. E. and Zajicek, J. P. (2012). The impact of motor and non motor symptoms on health state values in newly diagnosed idiopathic parkinson's disease, *Journal of Neurology* **259**(3): 462–468.

Shojaie, A. and Michailidis, G. (2010). Penalized likelihood methods for estimation of sparse high-dimensional directed acyclic graphs, *Biometrika* **97**(3): 519–538.

Tibar, H., El Bayad, K., Bouhouche, A., Ait Ben Haddou, E. H., Benomar, A., Yahyaoui, M., Benazzouz, A. and Regragui, W. (2018). Non-motor symptoms of parkinson's disease and their impact on quality of life in a cohort of moroccan patients, *Frontiers in neurology* **9**: 170–170.

van de Geer, S. and Bühlmann, P. (2013). $\ell_0$-penalized maximum likelihood for sparse directed acyclic graphs, *The Annals of Statistics* **41**(2): 536–567.

Wang, Y., Segarra, S. and Uhler, C. (2020). High-dimensional joint estimation of multiple directed gaussian graphical models, *Electronic Journal of Statistics* **14**(1): 2439–2483.

Wei, L., Hu, X., Zhu, Y., Yuan, Y., Liu, W. and Chen, H. (2017). Aberrant intra-and internetwork functional connectivity in depressed Parkinson's disease, *Scientific reports* **7**(1): 1–12.

Yang, Y., Wainwright, M. J. and Jordan, M. I. (2016). On the computational complexity of high-dimensional bayesian variable selection, *The Annals of Statistics* **44**(6): 2497–2532.

Yu, G. and Bien, J. (2017). Learning local dependence in ordered data, *Journal of Machine Learning Research* **18**(42): 1–60.

Zang, Y., Jiang, T., Lu, Y., He, Y. and Tian, L. (2004). Regional homogeneity approach to fmri data analysis, *NeuroImage* **22**(1): 394 – 400.

Zellner, A. (1986). On assessing prior distributions and bayesian regression analysis with g-prior distributions, *Bayesian inference and decision techniques: Essays in Honor of Bruno De Finetti* **6**: 233–243.