

Consistent Bayesian Sparsity Selection for High-dimensional Gaussian DAG Models with Multiplicative and Beta-mixture Priors

Xuan Cao

Kshitij Khare

Malay Ghosh

University of Cincinnati

University of Florida

University of Florida

March 11, 2019

Abstract

Estimation of the covariance matrix for high-dimensional multivariate datasets is a challenging and important problem in modern statistics. In this paper, we focus on high-dimensional Gaussian DAG models where sparsity is induced on the Cholesky factor L of the inverse covariance matrix. In recent work, ([Cao, Khare, and Ghosh, 2019]), we established high-dimensional sparsity selection consistency for a hierarchical Bayesian DAG model, where an Erdos-Renyi prior is placed on the sparsity pattern in the Cholesky factor L , and a DAG-Wishart prior is placed on the resulting non-zero Cholesky entries. In this paper we significantly improve and extend this work, by (a) considering more diverse and effective priors on the sparsity pattern in L , namely the beta-mixture prior and the multiplicative prior, and (b) establishing sparsity selection consistency under significantly relaxed conditions on p , and the sparsity pattern of the true model. We demonstrate the validity of our theoretical results via numerical simulations, and also use further simulations to demonstrate that our sparsity selection approach is competitive with existing state-of-the-art methods including both frequentist and Bayesian approaches in various settings.

1 Introduction

Covariance estimation and selection is a fundamental problem in multivariate statistical inference, and plays a crucial role in many data analytic methods. In high-dimensional settings, where the number of variables is much larger than the number of samples, the sample covariance matrix (traditional estima-

tor for the population covariance matrix) can perform rather poorly. See [Bickel and Levina, 2008a,b, El Karoui, 2007] for example. To address the challenge posed by high-dimensionality, several promising methods have been proposed in the literature. In particular, methods inducing sparsity in the covariance matrix Σ , its inverse Ω , or the Cholesky factor of the inverse, have proven to be very effective in applications. In this paper, we focus on imposing sparsity on the Cholesky factor of the inverse covariance (precision) matrix. These models are also referred to as Gaussian DAG models.

Consider a case when we have i.i.d. observations $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n$ obeying a p -variate normal distribution with mean vector $\mathbf{0}$ and covariance matrix Σ . Let $\Omega = LD^{-1}L^T$ be the unique modified Cholesky decomposition of the inverse covariance matrix $\Omega = \Sigma^{-1}$, where L is a lower triangular matrix with unit diagonal entries, and D is a diagonal matrix with positive diagonal entries. A given sparsity pattern on L corresponds to certain conditional independence relationships, which can be encoded in terms of a directed acyclic graph \mathcal{D} on the set of p variables as follows: if the variables i and j do not share an edge in \mathcal{D} , then $L_{ij} = 0$ (see Section 2 for more details).

There are two major approaches in the literature for sparse estimation of L . The first approach is based on regularized likelihood/pseudolikelihood using ℓ_1 penalization. See [Huang, Liu, Pourahmadi, and Liu, 2006, Rutimann and Buhlmann, 2009, Shojaie and Michailidis, 2010, Rothman, Levina, and Zhu, 2010, Aragam, Amini, and Zhou, 2015, Yu and Bien, 2016, Khare, Oh, Rahman, and Rajaratnam, 2017]. Some of these frequentist approaches assume L is banded, i.e., the elements of L that are far from the diagonal are taken to be zero. The other methods put restrictions on the maximum number of non-zero entries in L .

On the Bayesian side, when the underlying graph is known, literature exists that explores the posterior convergence rates for Gaussian concentration graph models, which induce sparsity in the inverse covariance matrix Ω . See [Banerjee and Ghosal, 2014, 2015, Xiang, Khare, and Ghosh, 2015, Lee and Lee, 2017] for example. Gaussian concentration graph models and Gaussian DAG models studied in this paper intersect only at perfect DAG models, which are equivalent to decomposable concentration graphical models. For general Gaussian DAG models, comparatively fewer works have tackled with asymptotic consistency properties. Recently, Cao, Khare, and Ghosh [2019] establish both strong model selection consistency and posterior convergence rates for sparse Gaussian DAG models in a high-dimensional regime. In particular, the authors consider a hierarchical Gaussian DAG model with DAG-Wishart priors introduced in [Ben-David, Li, Massam, and Rajaratnam, 2016] on the Cholesky parameter space and independent Bernoulli(q) priors for each edge in the DAG (the so-called Erdos-Renyi prior). However, the sparsity assumptions on the true model required to establish consistency are rather restrictive. In addition, as

a result of the extremely small value of the edge probability q in the Bernoulli prior, the simulation studies always tend to favor more sparse models under smaller values of p . Lee, Lee, and Lin [2018] also explore the Cholesky factor selection consistency under the empirical sparse Cholesky (ESC) prior and α -posteriors. Compared with [Cao, Khare, and Ghosh, 2019], under relaxed conditions in terms of the dimensionality, sparsity and lower bound of the non-zero elements in the Cholesky factor, Lee, Lee, and Lin [2018] establish strong model selection consistency with the α -posterior distribution.

It recently came to our attention that two more flexible alternative priors compared to the Erdos-Renyi prior have been considered in the undirected graphical models literature: (a) the multiplicative prior [Tan, Jasra, De Iorio, and Ebbels, 2017], and (b) the beta-mixture prior [Carvalho and Scott, 2009]. Both priors are more diverse than the Erdos-Renyi prior (the Erdos-Renyi prior can be obtained as a degenerate version of these priors), and have various attractive properties. For example, the multiplicative model prior can account for greater variability in the degree distribution as compared to the Erdos-Renyi model, while the beta-mixture prior allows for stronger control over the number of spurious edges and corrects for multiple hypothesis testing automatically. We provide the algebraic forms of these priors in Section 3 and Section 5 respectively, and refer the reader to [Carvalho and Scott, 2009, Tan, Jasra, De Iorio, and Ebbels, 2017] for a detailed discussion of their properties.

To the best of our knowledge, a rigorous investigation of high-dimensional posterior consistency properties with the multiplicative prior or the beta-mixture prior has not been undertaken for either undirected graphical models or Gaussian DAG models. Hence, our goal was to investigate if high-dimensional consistency results could be established under these two more diverse and algebraically complex class of prior distributions in the Gaussian DAG model setting. Another goal was to investigate if these high-dimensional posterior consistency results can be obtained under much weaker conditions as compared to [Cao, Khare, and Ghosh, 2019], particularly conditions similar to those in [Lee, Lee, and Lin, 2018]. This was a challenging goal, particularly for the multiplicative model prior, as the prior mass function is not available in closed form (note that the mass functions for the Erdos-Renyi, ESC and beta-mixture priors are available in closed form).

As the main contributions of this paper, we establish high-dimensional posterior consistency results for Gaussian DAG models with spike and slab priors on the Cholesky factor L , under both the multiplicative prior as well as the beta-mixture prior on the sparsity pattern in L (Theorems 4.1 to 5.3), using assumptions similar to those in [Lee, Lee, and Lin, 2018] (where a different setting of ESC priors and α -posteriors is used). Also, through simulation studies, we demonstrate that the models studied in this paper can outperform existing state-of-the-art methods including both penalized likelihood and Bayesian

approaches in different settings.

The rest of paper is organized as follows. Section 2 provides background material regarding Gaussian DAG model and introduce the spike and slab prior on the Choleksy factor. In Section 3, we revisit the multiplicative prior, and present our hierarchical Bayesian model and the parameter class for the inverse covariance matrices. Model selection consistency results for both the multiplicative prior and the beta-mixture prior are stated in Section 4 and Section 5 with proofs provided in Section 7. In Section 6 we use simulation experiments to illustrate the posterior ratio consistency result, and demonstrate the benefits of our Bayesian approach and computation procedures for Choleksy factor selection vis-a-vis existing Bayesian and penalized likelihood approaches. We end our paper with a discussion session in Section 8.

2 Preliminaries

In this section, we provide the necessary background material from graph theory, Gaussian DAG models, and also introduce our spike and slab prior on the Cholesky parameter.

2.1 Gaussian DAG Models

We consider the multivariate Gaussian distribution

$$\mathbf{Y} \sim N_p(0, \Omega^{-1}), \tag{1}$$

where Ω is a $p \times p$ inverse covariance matrix. Any positive definite matrix Ω can be uniquely decomposed as $\Omega = LD^{-1}L^T$, where L is a lower triangular matrix with unit diagonal entries, and D is a diagonal matrix with positive diagonal entries. This decomposition is known as the modified Cholesky decomposition of Ω (see for example Pourahmadi [2007]). In particular, the model (1) can be interpreted as a Gaussian DAG model depending on the sparsity pattern of L .

A directed acyclic graph (DAG) $\mathcal{D} = (V, E)$ consists of the vertex set $V = \{1, \dots, p\}$ and an edge set E such that there is no directed path starting and ending at the same vertex. As in [Ben-David, Li, Massam, and Rajaratnam, 2016, Cao, Khare, and Ghosh, 2019], we will without loss of generality assume a parent ordering, where that all the edges are directed from larger vertices to smaller vertices. For several applications in genetics, finance, and climate sciences, a location or time based ordering of variables is naturally available. For example, in genetic datasets, the variables can be genes or SNPs located contiguously on a chromosome, and their spatial location provides a natural ordering. More

examples can be found in [Huang, Liu, Pourahmadi, and Liu, 2006, Shojaie and Michailidis, 2010, Yu and Bien, 2016, Khare, Oh, Rahman, and Rajaratnam, 2017]. The set of parents of i , denoted by $pa_i(\mathcal{D})$, is the collection of all vertices which are larger than i and share an edge with i . Similarly, the set of children of i , denoted by $chi_i(\mathcal{D})$, is the collection of all vertices which are smaller than i and share an edge with i .

A Gaussian DAG model over a given DAG \mathcal{D} , denoted by $\mathcal{N}_{\mathcal{D}}$, consists of all multivariate Gaussian distributions which obey the directed Markov property with respect to a DAG \mathcal{D} . In particular, if $\mathbf{Y} = (Y_1, \dots, Y_p)^T \sim N_p(0, \Sigma)$ and $N_p(0, \Sigma = \Omega^{-1}) \in \mathcal{N}_{\mathcal{D}}$, then

$$Y_i \perp \mathbf{Y}_{\{i+1, \dots, p\} \setminus pa_i(\mathcal{D})} | \mathbf{Y}_{pa_i(\mathcal{D})},$$

for each $1 \leq i \leq p$. Furthermore, it is well-known that if $\Omega = LD^{-1}L^T$ is the modified Cholesky decomposition of Ω , then $N_p(0, \Omega^{-1}) \in \mathcal{N}_{\mathcal{D}}$ if and only if $L_{ij} = 0$ whenever $i \notin pa_j(\mathcal{D})$. In other words, the structure of the DAG \mathcal{D} is uniquely reflected in the sparsity pattern of the Cholesky factor L . In light of this, it is often more convenient to reparametrize the inverse covariance matrix in terms of the Cholesky parameter (L, D) .

2.2 Notations

Consider the modified cholesky decomposition $\Omega = LD^{-1}L^T$, where L is a lower triangular matrix with all the unit diagonals and $D = \text{Diag} \{d_1, d_2, \dots, d_p\}$, where d_i 's are all positive. We suggest to impose spike and slab priors on the lower diagonal of L to recover the sparse structure of the Cholesky factor. To facilitate this purpose, we introduce latent binary variables $Z = \{Z_{21}, \dots, Z_{kj}, \dots, Z_{p,p-1}\}$ for $1 \leq j < k \leq p$ to indicate whether L_{kj} is active, i.e., $Z_{kj} = 1$ if $L_{kj} \neq 0$ and 0, otherwise. We can view the binary variable Z_{kj} as the indicator for the sparsity pattern of L . In other words, for each $1 \leq j \leq p-1$, let Z_j , a subset of $\{j+1, j+2, \dots, p\}$, be the index set of all non-zero components in $\{Z_{j+1,j}, \dots, Z_{p,j}\}$. Z_j explicitly gives the support of the Cholesky factor and the sparsity pattern of the underlying DAG. Denote $|Z_j| = \sum_{k=j+1}^p Z_{kj}$ as the cardinality of set Z_j for $1 \leq j \leq p-1$.

Following the definition of Z , for any $p \times p$ matrix A , denote the column vectors $A_{Z,j}^> = (A_{kj})_{k \in Z_j}$ and $A_{Z,i}^> = (A_{ii}, (A_{Z,i}^>)^T)^T$. Also, let $A_Z^{>j} = (A_{ki})_{k,i \in Z_j}$,

$$A_Z^{>i} = \begin{bmatrix} A_{ii} & (A_{Z,i}^>)^T \\ A_{Z,i}^> & A_Z^{>i} \end{bmatrix}.$$

In particular, $A_{\bar{Z},p}^{\geq} = A_{\bar{Z}}^{\geq q} = A_{pp}$.

Next, we provide some additional required notation. For $x \in \mathbb{R}^p$, let $\|x\|_r = \left(\sum_{j=1}^p |x_j|^r\right)^{\frac{1}{r}}$ and $\|x\|_{\infty} = \max_j |x_j|$ represent the standard l_r and l_{∞} norms. For a $p \times p$ matrix A , let $eig_1(A) \leq eig_2(A) \dots eig_p(A)$ be the ordered eigenvalues of A and denote

$$\|A\|_{\max} = \max_{1 \leq i, j \leq p} |A_{ij}|,$$

$$\|A\|_{(r,s)} = \sup \{\|Ax\|_s : \|x\| = 1\}, \text{ for } 1 \leq r, s < \infty.$$

In particular,

$$\|A\|_{(1,1)} = \max_j \sum_i |A_{ij}|, \quad \|A\|_{(\infty,\infty)} = \max_i \sum_j |A_{ij}| \text{ and } \|A\|_{(2,2)} = eig_p(A)^{\frac{1}{2}}.$$

2.3 Spike and Slab Prior on Cholesky Parameter

In this section, we specify our spike and slab prior on the Cholesky factor as follows.

$$L_{kj} \mid d_j, Z_{kj} \stackrel{ind}{\sim} Z_{kj} N(0, \tau^2 d_j) + (1 - Z_{kj}) \delta_0(L_{kj}), \quad 1 \leq j < k \leq p, \quad (2)$$

$$d_j \stackrel{ind}{\sim} \text{Inverse-Gamma}(\lambda_1, \lambda_2), \quad j = 1, 2, \dots, p, \quad (3)$$

for some constants $\tau, \lambda_1, \lambda_2 \geq 0$, where $\delta_0(L_{kj})$ denotes a point mass at 0. We refer to (2) and (3) as our spike and slab Cholesky (SSC) prior. $Z_{jk} = 1$ implies L_{jk} being the ‘‘signal’’ (i.e., from the slab component), and $Z_{jk} = 0$ implies L_{jk} being the noise (i.e., from the spike component). Note that to obtain our desired asymptotic consistency results, appropriate conditions for these hyperparameters will be introduced in Section 4.1. Xu and Ghosh [2015] also impose this type of priors on the regression factors. Further comparisons and discussion are provided in Remark 3.

Remark 1. *Note that in (3), we are allowing the hyperparameters for the inverse-gamma prior to be zero. In [Cao, Khare, and Ghosh, 2019], the DAG-Wishart prior with multiple shape parameters introduced in [Ben-David, Li, Massam, and Rajaratnam, 2016] is placed on the Cholesky parameter. As indicated in Theorem 7.3 in [Ben-David, Li, Massam, and Rajaratnam, 2016], the DAG-Wishart distribution defined on the Cholesky parameter space given a DAG yields the independent inverse-gamma distribution with strictly positive shape and scale parameters on d_j and multivariate Gaussian distribution on the non-zero*

elements in each column of L given d_j . Hence, for given DAG structures, there are some difference and connection between the DAG-Wishart prior and our spike and slab prior.

3 Model Specification

In this section, we revisit the multiplicative prior introduced in [Tan, Jasra, De Iorio, and Ebbels, 2017] over space of graphs, and specify our hierarchical model.

3.1 Multiplicative Prior

In the context of Gaussian graphical model, Tan, Jasra, De Iorio, and Ebbels [2017] allow the probability of a link between nodes k and j , q_{kj} to vary with i, j by taking $q_{kj} = \omega_k \omega_j$ and $0 < \omega_j < 1$ for each $1 \leq j \leq p$. The authors further treat each ω_i as a variable with a beta prior to adopt a fully Bayesian approach. The authors further utilize Laplace approximations, and through simulation studies, show that the proposed multiplicative model (following the nomenclature in [Tan, Jasra, De Iorio, and Ebbels, 2017]) facilitates the purpose to encourage sparsity or graphs that exhibit particular degree patterns based on prior knowledge. Adapted to our framework, we consider the following multiplicative prior over the space of sparsity variation for the Cholesky factor.

$$\pi(Z \mid \omega_1, \dots, \omega_p) = \prod_{1 \leq j < k \leq p} (\omega_k \omega_j)^{Z_{kj}} (1 - \omega_k \omega_j)^{1 - Z_{kj}}, \quad (4)$$

$$\omega_j \sim \text{Beta}(\alpha_1, \alpha_2), \quad 1 \leq j \leq p, \quad (5)$$

where α_1, α_2 are positive constants. Compared with the universal indicator probability q in an Erdos-Renyi prior, here we allow the variation attainable in the degree structure of each node through different values of ω_j . Note that under the multiplicative prior, the marginal posterior for Z can not be obtained in closed form, which leads to further challenges not only in the theoretical analysis, but also in the computational strategy. We will elaborate on this matter in Section 6.

3.2 Hierarchical Model Formulation

Let $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n$ be independent and identically distributed p -variate Gaussian vectors with mean 0 and true covariance matrix $\Sigma_0 = (\Omega_0)^{-1}$, where $\Omega_0 = L_0(D_0)^{-1}(L_0)^T$ is the modified Cholesky decomposition of Ω_0 . Let $S = \frac{1}{n} \sum_{i=1}^n \mathbf{Y}_i \mathbf{Y}_i^T$ denotes the sample covariance matrix. The sparsity pattern of the true Cholesky factor L_0 is uniquely encoded in the true binary variable denoted as Z_0 . Similar to [Cao, Khare,

and Ghosh, 2019], we also denote d as the maximum number of non-zero entries in any column of L_0 , and $s = \min_{1 \leq j, i \leq p, i \in Z_j} |(L_0)_{ji}|$. For sequences a_n and b_n , $a_n \sim b_n$ means $\frac{a_n}{b_n} \rightarrow c$ for some constant $c > 0$, as $n \rightarrow \infty$. Let $a_n = o(b_n)$ represent $\frac{a_n}{b_n} \rightarrow 0$ as $n \rightarrow \infty$.

The class of spike and slab Cholesky distributions in Section 2 and the multiplicative priors in Section 3.1 can be used for Bayesian model selection of the Cholesky factor through the following hierarchical model,

$$\mathbf{Y} \mid (D, L), Z \sim N_p(\mathbf{0}, (LD^{-1}L^T)^{-1}), \quad (6)$$

$$L_{kj} \mid d_j, Z_{kj} \stackrel{\text{ind}}{\sim} Z_{kj}N(\mathbf{0}, \tau^2 d_j) + (1 - Z_{kj})\delta_0(L_{kj}), \quad 1 \leq j < k \leq p, \quad (7)$$

$$d_j \stackrel{\text{ind}}{\sim} \text{Inverse-Gamma}(\lambda_1, \lambda_2), \quad j = 1, 2, \dots, p, \quad (8)$$

$$\pi(Z \mid \omega_1, \dots, \omega_p) = \prod_{1 \leq j < k \leq p} (\omega_k \omega_j)^{Z_{kj}} (1 - \omega_k \omega_j)^{1 - Z_{kj}}, \quad (9)$$

$$\omega_j \sim \text{Beta}(\alpha_1, \alpha_2), \quad 1 \leq j \leq p, \quad (10)$$

where $\text{Beta}(\alpha_1, \alpha_2)$ represents the beta distribution with shape parameters α_1, α_2 . The proposed hierarchical model now has five hyperparameters: the scale parameter $\tau > 0$ in model (7) controlling the variance of the spike part in the spike and slab prior on each L_{kj} , the shape parameter λ_1 and scale parameter λ_2 in model (8), and the two positive shape parameters in the beta distribution in model (10). Further restrictions on these hyperparameters to ensure desired consistency will be specified in Section 4.1.2.

The intuition behind this set-up with latent variables is that the elements in the Cholesky factor L with zero or very small values will be identified with zero Z values, while the active entries will be classified as $Z = 1$. We use the posterior probabilities of all the $\frac{p(p-1)}{2}$ latent variables Z to identify the active elements in L . In particular, the following lemmas help specify the upper bound for the marginal probability ratio and the marginal posterior ratio for any “non-true” model Z compared with the true model Z under the multiplicative prior. The proof will be provided in Section 7.1.

Lemma 3.1. *If the hyperparameter α_2 in model (10) satisfies $\alpha_2 \sim \max\{p^c, d^{\frac{2c}{c-2}}\}$, for $c > 2$, we have*

$$\frac{\pi(Z)}{\pi(Z_0)} \leq e^{2\alpha_1^2 + 2\alpha_1 + \frac{2}{\alpha_1}} \prod_{j=1}^p \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)}, \quad (11)$$

for $p \geq 4 + \frac{4}{\alpha_1} + 2\sqrt{\alpha_1}$.

Lemma 3.1 further enables the marginalized posterior likelihood ratio to be upper bounded by decomposed

prior terms absorbed into the product of items as follows.

Lemma 3.2. *If $\alpha_2 \sim \max\{p^c, d^{\frac{2c}{c-2}}\}$, for $c > 2$, the marginal posterior ratio between any “non-true” model Z and the true model Z_0 under the multiplicative prior in (4) and (5) satisfies*

$$\begin{aligned}
& \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\
& \leq M_1 \prod_{j=1}^{p-1} (n\tau^2)^{-\frac{|Z_j| - |Z_{0j}|}{2}} \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \\
& \quad \times \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \left(\frac{\tilde{S}_{j|Z_{0j}}}{\tilde{S}_{j|Z_j}} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_{j|Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}}{\tilde{S}_{j|Z_j} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}} \right)^{\frac{n}{2} + \lambda_1} \\
& \triangleq M_1 \times \prod_{j=1}^{p-1} PR'_j(Z, Z_0), \tag{12}
\end{aligned}$$

where $M_1 = e^{2\alpha_1^2 + 2\alpha_1 + \frac{2}{\alpha_1}}$, $\tilde{S} = S + \frac{1}{n\tau_{n,p}^2} I_p$ and $\tilde{S}_{j|Z_j} = \tilde{S}_{jj} - (\tilde{S}_{Z \cdot j}^{\geq})^T (\tilde{S}_Z^{\geq j})^{-1} \tilde{S}_{Z \cdot j}^{\geq}$.

4 Model Selection Consistency

In this section we will explore the high-dimensional asymptotic properties of the Bayesian model selection approach for the Cholesky factor specified in Section 3.2. For this purpose, we will work in a setting where the dimension $p = p_n$ of the data vectors, and the hyperparameters vary with the sample size n and $p_n \geq n$. Assume that the data is actually being generated from a true model specified as follows. Let $\mathbf{Y}_1^n, \mathbf{Y}_2^n, \dots, \mathbf{Y}_n^n$ be independent and identically distributed p_n -variate Gaussian vectors with mean 0 and true covariance matrix $\Sigma_0^n = (\Omega_0^n)^{-1}$, where $\Omega_0^n = L_0^n (D_0^n)^{-1} (L_0^n)^T$ is the modified Cholesky decomposition of Ω_0^n . The sparsity pattern of the true Cholesky factor L_0^n is reflected in Z_0^n . Recall the definition in Section 3.2 that d_n is the maximum number of non-zero entries in any column of L_0^n , and $s_n = \min_{1 \leq j, i \leq p, i \in Z_j} |(L_0^n)_{ji}|$. In order to establish our asymptotic consistency results, we need the following mild assumptions with respective discussion/interpretation.

4.1 Assumptions

4.1.1 Assumptions on the True Parameter Class

Assumption 1. *There exists $\epsilon_0 \leq 1$, such that for every $n \geq 1$, $0 < \epsilon_0 \leq \text{eig}_1(\Omega_0^n) \leq \text{eig}_{p_n}(\Omega_0^n) \leq \epsilon_0^{-1}$.*

This assumption ensures that the eigenvalues of the true precision matrices are bounded by fixed constants, which has been commonly used for establish high dimensional covariance asymptotic properties.

See for example [Bickel and Levina, 2008a, El Karoui, 2008, Banerjee and Ghosal, 2014, Xiang, Khare, and Ghosh, 2015, Banerjee and Ghosal, 2015]. Previous work [Cao, Khare, and Ghosh, 2019] relaxes this assumption by allowing the lower and upper bounds on the eigenvalues to depend on p and n .

Assumption 2. $d_n \sqrt{\frac{\log p_n}{n}} \rightarrow 0$ as $n \rightarrow \infty$.

This is a much weaker assumption for high dimensional covariance asymptotic than for example, [Xiang, Khare, and Ghosh, 2015, Banerjee and Ghosal, 2014, 2015, Cao, Khare, and Ghosh, 2019]. Here we essentially allow the number of variables p_n to grow slower than e^{n/d_n^2} compared to previous literatures with rate e^{n/d_n^4} .

Assumption 3. $\frac{d_n \log p_n}{s_n^2 n} \rightarrow 0$ as $n \rightarrow \infty$.

Recall that s_n is the smallest (in absolute value) non-zero off-diagonal entry in L_0^n . Hence, this assumption also known as the “beta-min” condition also provides a lower bound for the “slab” part of L_0^n that is needed for establishing consistency. This type of condition has been used for the exact support recovery of the high-dimensional linear regression models as well as Gaussian DAG models. See for example [Lee, Lee, and Lin, 2018, Yang, Wainwright, and Jordan, 2016, Khare, Oh, Rahman, and Rajaratnam, 2017, Cao, Khare, and Ghosh, 2019, Yu and Bien, 2016].

Remark 2. *It is worthwhile to point out that our assumptions on the true Cholesky factor are weaker compared to [Lee, Lee, and Lin, 2018]. In particular, Lee, Lee, and Lin [2018] introduce conditions A(2) and A(4) on the sparsity pattern of the true Cholesky factor such that the number of non-zero elements in each row as well as each column of L_0^n to be smaller than some constant s_0 , while in this paper, we are allowing the maximum number of non-zero entries in any column of L_0^n to grow at a smaller rate than $\sqrt{\frac{n}{\log p_n}}$ (Assumption 2).*

4.1.2 Assumptions on the Prior Hyperparameters

Assumption 4. $\pi(Z) = 0$ for all Z satisfying $\max_{1 \leq j \leq p-1} |Z_j| \geq R_n$, where $R_n \sim n(\log n)^{-1}$.

This assumption essentially states that the prior on the space of the $2^{\binom{p_n}{2}}$ possible models, places zero mass on unrealistically large models (see similar assumptions in [Johnson and Rossell, 2012, Shin, Bhattacharya, and Johnson, 2018, Narisetty and He, 2014] in the context of regression). Assumption 4 is also more relaxed compared with Condition (P) in [Lee, Lee, and Lin, 2018] where $R_n \sim n(\log p)^{-1}\{(\log n)^{-1} \vee c_3\}$ for some constant c_3 . Note that this condition is for the hyperparameter of the prior distribution on the latent variables only, which does not affect the true parameter space.

Assumption 5. The hyperparameter τ_{n,p_n} in (9) satisfies $\frac{d_n}{\tau_{n,p_n}^2 \log p_n} \rightarrow 0$ and $\frac{\sqrt{\frac{n}{\tau_{n,p_n}^2}}}{p_n^{\frac{(1-1/\kappa)c}{2}} \log n} \rightarrow 0$, as $n \rightarrow \infty$, for some constant $\kappa > 1$.

This assumption provides the rate at which the variance of the slab prior is required to grow to guarantee desired model selection consistency. Similar conditions on the hyperparameter can be seen in [Narisetty and He, 2014, Shin, Bhattacharya, and Johnson, 2018, Johnson and Rossell, 2012].

Assumption 6. There exists a constant $c > 0$, such that the hyperparameters in model (8) satisfy $0 \leq \lambda_{1n}, \lambda_{2n} < c$ and the shape parameters in model (10) satisfies $0 < \alpha_{1n} < c$, $\alpha_2 \sim \max \left\{ p_n^c, d_n^{\frac{2c}{c-2}} \right\}$, for $c > 2\kappa$, for some $\kappa > 1$.

This assumption provides the rate at which the shape parameter needs to grow to ensure desired consistency. Previous literature with Erdos-Renyi priors puts restrictions on the rate of the edge probability. In particular, previous work [Cao, Khare, and Ghosh, 2019] assumes $q = e^{-\eta_n n}$, where $\eta_n = d_n \left(\frac{\log p_n}{n} \right)^{\frac{1/2}{1+k/2}}$ for some $k > 0$ to penalize large models. Similar assumptions on the hyperparameters can be also found in [Yang, Wainwright, and Jordan, 2016, Narisetty and He, 2014] under regression setting. In Section 6.2, we will see the proposed model without specifying particular values for q helps avoiding the potential computation limitation such as simulation results always favor the most sparse model.

For the rest of this paper, $p_n, \Omega_0^n, \Sigma_0^n, L_0^n, D_0^n, Z_0^n, Z^n, d_n, \tau_n, s_n, \alpha_{1n}, \alpha_{2n}$ will be denoted as $p, \Omega_0, \Sigma_0, L_0, D_0, Z_0, Z, d, \tau, s, \alpha_1, \alpha_2$ by leaving out the superscript for notational convenience.

4.2 Posterior Ratio Consistency

We now state and prove the main model selection consistency results. The proofs for all the theorems will be provided in Section 7.1 and Section 7.2. Our first result establishes what we refer to as ‘‘posterior ratio consistency’’ (following the terminology in [Cao, Khare, and Ghosh, 2019]). This notion of consistency implies that the true model will be the mode of the posterior distribution among all the models with probability tending to 1 as $n \rightarrow \infty$.

Theorem 4.1. Under Assumptions 1-6, the following holds:

$$\max_{Z \neq Z_0} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \xrightarrow{\bar{P}} 0, \text{ as } n \rightarrow \infty.$$

4.3 Model Selection Consistency for Posterior Mode

If one was interested in a point estimate of Z which reflects the sparsity pattern of L_0 , the most apparent choice would be the posterior mode defined as

$$\hat{Z} = \arg \max_Z \pi(Z|\mathbf{Y}). \quad (13)$$

From a frequentist point of view, it would be natural to obtain if we have model selection consistency for the posterior mode, which follows immediately from posterior ratio consistency established in Theorem 4.1, by noting that $\max_{Z \neq Z_0} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} < 1 \Rightarrow \hat{Z} = Z_0$. Therefore, we have the following corollary.

Corollary 4.1. *Under Assumptions 1-6, the posterior mode \hat{Z} is equal to the true model Z_0 with probability tending to 1, i.e.,*

$$\bar{P}(\hat{Z} = Z_0) \rightarrow 1, \text{ as } n \rightarrow \infty.$$

Remark 3. *In the context of linear regression, Xu and Ghosh [2015] tackle the Bayesian group lasso problem. In particular, the authors propose the following hierarchical Bayesian model:*

$$\begin{aligned} \mathbf{Y} | X, \beta, \sigma^2 &\sim N(X\beta, \sigma^2 I) \\ \beta_g &\stackrel{\text{ind}}{\sim} \sigma^2, \tau_g^2 \sim (1 - \pi_0)N(0, \sigma^2 \tau_g^2 I) + \pi_0 \delta_0(\beta_g), \quad g = 1, 2, \dots, G, \\ \tau_g^2 &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{m_g + 1}{2}, \frac{\lambda^2}{2}\right), \quad g = 1, 2, \dots, G, \\ \sigma^2 &\stackrel{\text{ind}}{\sim} \text{Inverse-Gamma}(\alpha_1, \alpha_2). \end{aligned}$$

In particular, they impose an independent spike and slab type prior on each factor β_g (conditional on the variance parameter σ^2), and an inverse Gamma prior on the variance. Each regression factor is explicitly present in the model with a probability π_0 . In this setting under an orthogonal design, the authors in [Xu and Ghosh, 2015] establish oracle property and variable selection consistency for the median thresholding estimator of the regression coefficients on the group level. Note that with parent ordering, the off-diagonal entries in the i^{th} column of L can be interpreted as the linear regression coefficients corresponding to fitting the i^{th} variable against all variables with label greater than i . Hence, there are similarities with respect to the model and consistency results between [Xu and Ghosh, 2015] and this work. However, despite these similarities, fundamental differences exist in these models and the corresponding analysis. Firstly, the number of groups (or factors) is considered to be fixed in [Xu and Ghosh, 2015], while we allow the number of predictors to grow at an exponential rate of n in a ultra high-dimensional setting, which creates

more theoretically challenges. Secondly, the ‘design’ matrices corresponding to the regression coefficients in each column of L which can be represented as functions of the sample covariance matrix S are random and correlated with each other, while [Xu and Ghosh, 2015] only considers the orthogonal design where $X^T X = I$ with no correlation introduced. Thirdly, the consistency result in [Xu and Ghosh, 2015] focuses only on group level selection only and is tailored for problems that only require group level sparsity, while our model can induce sparsity in each individual element of L . The authors also propose a Bayesian hierarchical model referred to as Bayesian sparse group lasso to enable shrinkage both at the group level and within a group. However, no consistency results are addressed regarding this model. Lastly, in our model, each coefficient is present independently with multiplicative prior that incorporates information that L is sparse, which is not the case in [Xu and Ghosh, 2015] as each factor is present with $\pi_0 = 0.5$. In particular, all the aspects discussed above lead to major differences and further challenges in analyzing the ratio of posterior probabilities.

4.4 Strong Model Selection Consistency

Next we establish another stronger result (compared to Theorem 4.1) which implies that the posterior mass assigned to the true model Z_0 converges to 1 in probability (under the true model). Following [Narisetty and He, 2014, Cao, Khare, and Ghosh, 2019], we refer to this notion of consistency as strong selection consistency.

Theorem 4.2. *Under Assumptions 1-6, the following holds:*

$$\pi(Z_0|\mathbf{Y}) \xrightarrow{P} 1, \text{ as } n \rightarrow \infty.$$

Remark 4. *We would like to point out that our posterior ratio consistency and strong model selection consistency do not require any additional assumptions on bounding the maximum number of edges. In particular, Cao, Khare, and Ghosh [2019] consider only the DAGs with the total number of edges at most $\frac{1}{8}d \left(\frac{n}{\log p}\right)^{\frac{1+k}{2+k}}$ for $k > 0$. By the assumptions in the previous work, it follows that the DAGs in the analysis do not include the models where the Cholesky factor has one or more non-zero elements for each column, since $p/\frac{1}{8}d \left(\frac{n}{\log p}\right)^{\frac{1+k}{2+k}} \rightarrow \infty$, as $n \rightarrow \infty$, while in our result, each row can have at most $R_n \sim \frac{n}{\log n}$ number of non-zero entries as indicated in Assumption 4. Hence, our strong model selection consistency results is more general than [Cao, Khare, and Ghosh, 2019, Lee, Lee, and Lin, 2018] in the sense that the consistency holds for a larger class of DAGs.*

5 Results for Beta-mixture Prior

Though the multiplicative prior could allow variation among the indicator probabilities, the intractable marginal posteriors remain problematic in practice. The authors in [Tan, Jasra, De Iorio, and Ebbels, 2017] address this issue via Laplace approximation. However, the computational cost for that will become extensive as p increases. To obtain the marginal posterior probabilities in closed form and for ease of computation, we consider the following beta-mixture prior over the space of Z introduced in [Carvalho and Scott, 2009],

$$Z_{kj} \mid q \stackrel{i.i.d.}{\sim} \text{Bern}(q), \quad 1 \leq j < k \leq p, \quad (14)$$

$$q \sim \text{Beta}(\alpha_1, \alpha_2), \quad (15)$$

where $\text{Bern}(q)$ denotes the Bernoulli distribution with probability q , and $\text{Beta}(\alpha_1, \alpha_2)$ represents the beta distribution with shape parameters α_1, α_2 . We refer to model (14) and (15) as the beta-mixture prior over the space of latent variables indicating the sparsity structure for the Cholesky factor.

Remark 5. *Cao, Khare, and Ghosh [2019], Banerjee and Ghosal [2015] introduce an Erdos-Renyi type of distribution on the space of DAGs as the prior distribution for DAGs, where each directed edge is present with probability q independently of the other edges. In particular, they define $\gamma_{ij} = \mathbb{I}\{(i, j) \in E(\mathcal{D})\}$, $1 \leq i < j \leq p$ to be the edge indicator and let γ_{ij} , $1 \leq i < j < p$ be independent identically distributed Bernoulli(q) random variables. Cao, Khare, and Ghosh [2019] establish the DAG selection consistency under suitable assumptions. while Banerjee and Ghosal [2015] address the estimation consistency, and provide high-dimensional Laplace approximations for the marginal posterior probabilities for the graphs. In our framework, we extend the previous work by putting a beta distribution on the edge probability q . The beta-mixture type of priors have previously been placed on graphs for simulation purpose in [Carvalho and Scott, 2009], but the theoretical properties have yet to be investigated. A clear advantage of such an approach as indicated in [Carvalho and Scott, 2009] is that treating the previous fixed tuning constant q as a model parameter shrinks the graph size to a data-determined value of q , and allows strong control over the number of spurious edges.*

In order to obtain the posterior consistency for Z , we need the following lemma, which specifies the closed form for the marginal posterior density of Z with proof provided in Section 7.3.

Lemma 5.1. *The marginal posterior density $\pi(Z|Y)$ under the beta-mixture prior satisfies*

$$\begin{aligned}
& \pi(Z|\mathbf{Y}) \\
& \propto B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right) \\
& \quad \times \prod_{j=1}^{p-1} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{p}{2}-\lambda_1} \frac{|\tilde{S}_Z^{>j}|^{-\frac{1}{2}}}{(n\tau^2)^{|Z_j|/2}} \\
& = B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^p |Z_j| \right) \\
& \quad \times \prod_{j=1}^{p-1} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{p}{2}-\lambda_1} \frac{\left(|\tilde{S}_Z^{>i}| \tilde{S}_{j|Z_j} \right)^{-\frac{1}{2}}}{(n\tau^2)^{|Z_j|/2}}, \tag{16}
\end{aligned}$$

in which $\tilde{S} = S + \frac{1}{n\tau^2} I_p$, $\tilde{S}_{j|Z_j} = \tilde{S}_{jj} - (\tilde{S}_{Z\cdot j}^{>})^T (\tilde{S}_{Z\cdot j}^{>j})^{-1} \tilde{S}_{Z\cdot j}^{>}$ and $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$. The second equation follows from $|\tilde{S}_Z^{>j}| = |\tilde{S}_Z^{>i}| \left(\tilde{S}_{jj} - (\tilde{S}_{Z\cdot j}^{>})^T (\tilde{S}_{Z\cdot j}^{>j})^{-1} \tilde{S}_{Z\cdot j}^{>} \right) = |\tilde{S}_Z^{>i}| \tilde{S}_{j|Z_j}$.

In particular, these posterior probabilities can be used to select a model representing the sparsity pattern of L by computing the posterior mode that maximize the posterior densities. The convenient closed form for the marginal posterior in (16) also yields nice posterior ratio consistency under the following weaker assumption on α_2 compared with Assumption 6.

Assumption 7. *There exists a constant $c > 0$, such that the hyperparameters in model (8) satisfy $0 \leq \lambda_{1n}, \lambda_{2n} < c$ and the shape parameters in model (10) satisfies $0 < \alpha_{1n} < c$, $\alpha_{2n} \sim p^c$.*

Theorem 5.2. *Under Assumptions 1-5 and 7, the following holds under the beta-mixture prior:*

$$\max_{Z \neq Z_0} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \xrightarrow{\bar{P}} 0, \text{ as } n \rightarrow \infty.$$

The next theorem establishes the strong selection consistency under the beta-mixture prior. See proofs for Theorem 5.2 and Theorem 5.3 in Section 7.3.

Theorem 5.3. *Under Assumptions 1-6, for the beta-mixture prior, the following holds:*

$$\pi(Z_0|\mathbf{Y}) \xrightarrow{\bar{P}} 1, \text{ as } n \rightarrow \infty.$$

Remark 6. *We would like to point out that posterior ratio consistency (Theorem 5.2 does not require any restriction on c (the rate of the shape parameter in the beta distribution (15)) that will be growing,*

this requirement is only needed for strong selection consistency (Theorem 5.3). Similar restrictions on the hyperparameters have been considered for establishing consistency properties in the regression setup. See [Yang, Wainwright, and Jordan, 2016, Lee, Lee, and Lin, 2018, Cao, Khare, and Ghosh, 2018] for example.

The closed form for the marginal posterior probability in (16) is convenient for showing the consistency. However, when it comes to simulation, the beta term in (16) pertaining to the beta-mixture prior is often too large, and could sometimes blow up when p is relatively large. In addition, for the beta-mixture prior, probability q is assumed to be universal across all indicators, which seems not flexible and diverse enough. In the following section, we will take on the task to investigate and evaluate the simulation performance for both the multiplicative model and the beta-mixture model.

6 Simulation Studies

In this section, we demonstrate our main results through simulation studies. First recall from (16) that the marginal posterior distributions for Z under the beta-mixture prior can be derived analytically in closed form (up to a constant) in (16). Therefore, we can evaluate the parameter space more clearly with this naturally assigned “score”, that is the posterior probability.

For the multiplicative prior, the ω_j ($1 \leq j \leq p$) can not be integrated out, thus the closed form for the marginal distribution of Z can not be conveniently acquired. As indicated in [Tan, Jasra, De Iorio, and Ebbels, 2017], evaluating the marginal densities via Monte Carlo becomes more computationally intensive as the dimension increases. Therefore, the authors propose to estimate these quantities efficiently through Laplace approximation instead. Detailed functional and Hessian expressions can be found in the supplemental material in [Tan, Jasra, De Iorio, and Ebbels, 2017]. Here we adopt the same Laplace approximation for estimating the marginal densities for Z . However, as we will see in Figure 3, though the multiplicative prior could potentially lead to better model selection performance, the additional procedure when evaluating each individual posterior probability could be quite time consuming. In particular, the Newton-type algorithm used for obtaining the mode of the log-likelihood runs extremely slow in higher dimensions.

6.1 Simulation I: Illustration of Posterior Ratio Consistency

In this section, we illustrate the consistency result in Theorem 4.1 and Theorem 5.2 using a simulation experiment. Our goal is to show that the log of the posterior ratio for any “non-true” model compared to

the true model will converge to negative infinity. To serve this purpose, we consider 10 different values of p ranging from 150 to 1500, and choose $n = p/3$. Next, for each fixed p , a $p \times p$ lower triangular matrix with diagonal entries 1 and off-diagonal entries 0.5 is constructed. In particular, unlike in previous work [Cao et al., 2019] where the expected value of non-zero entries in each column of L_0 does not exceed 3, here we randomly chose 3% or 5% of the lower triangular entries of the Cholesky factor and set them to be 0.5. The remaining entries were set to zero.

The purpose of this setting is to show our consistency requires more relaxed sparsity assumptions on the true model compared to [Cao, Khare, and Ghosh, 2019]. We refer to this matrix as L_0 . The matrix L_0 also reflects the true underlying DAG structure encoded in Z_0 . Next, we generate n i.i.d. observations from the $N(0_p, (L_0^{-1})^T L_0^{-1})$ distribution, and set the hyperparameters as $c = 2$, $\tau_{n,p} = \sqrt{n}$, $\lambda_1 = \lambda_2 = 0.05$, $\alpha_1 = 0.05$ for $i = 1, 2, \dots, p$. The above process ensures all the assumptions are satisfied. We then examine posterior ratio consistency under four different cases by computing the log posterior ratio of a “non-true” model Z and Z_0 as follows.

1. Case 1: Model Z is a submodel of Z_0 and the number of total non-zero entries of Z is exactly half of Z_0 , i.e. $\sum Z = \frac{1}{2} \sum Z_0$.
2. Case 2: Z_0 is a submodel of Z and the number of total non-zero entries of Z is exactly twice of Z_0 , i.e. $\sum Z = 2 \sum Z_0$.
3. Case 3: Z is not necessarily a submodel of Z_0 , but satisfying the number of total non-zero entries in Z is half the number of non-zero entries in Z_0 .
4. Case 4: Z_0 is not necessarily a submodel of Z , but the number of total non-zero entries in Z is twice the number of non-zero elements in Z_0 .

The log of the posterior probability ratio for various cases under two different sparsity settings and our two different priors are provided in Figure 1. As expected the log of the posterior probability ratio decreases to large negative numbers as n becomes large in all four cases and in both sparsity settings and under both sparsity priors, thereby providing a numerical illustration of Theorem 4.1.

We would like to point out that in [Cao, Khare, and Ghosh, 2019], the log of posterior ratios are almost all positive real numbers, when $p \leq 1500$ and the expected value of non-zero entries in each column of L_0 does not exceed 3, which indicates the hierarchical model with DAG-Wishart distribution and the Erdos-Renyi type of prior over graphs only performs better with really higher dimension and much more sparse settings. In particular, this leads to one potential drawback of using the DAG-Wishart

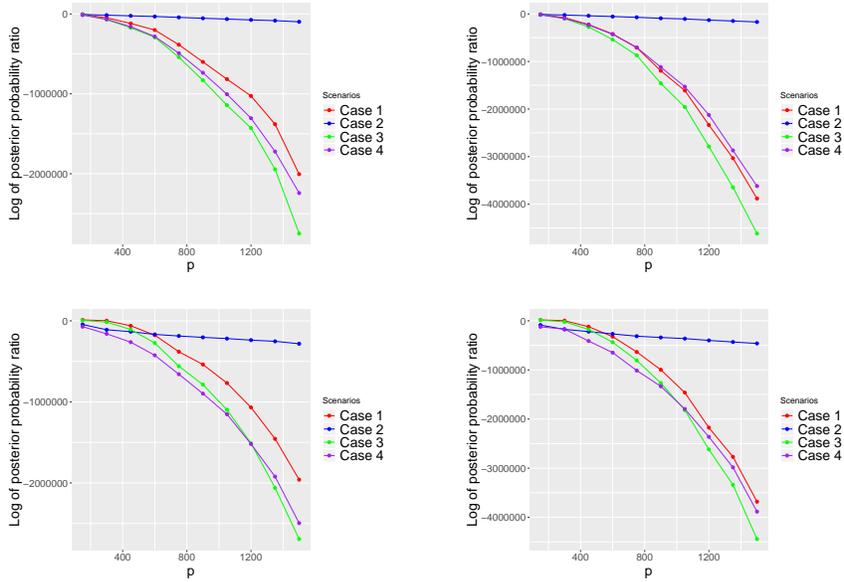


Figure 1: Log of posterior probability ratio for Z and Z_0 for various choices of the “non-true” model Z . Here Z_0 denotes the true underlying model indicator. Left: 3% sparsity; right: 5% sparsity; top: beta-mixture prior; bottom: multiplicative prior.

distribution coupled with the Erdos-Renyi type of prior on the Cholesky space, as in real applications, extremely high-dimensional and sparse data sets are not very commonly seen, while our spike and slab Cholesky prior with the beta-mixture or multiplicative prior is more adaptable and diverse in that aspect.

6.2 Simulation II: Illustration of Model Selection

In this section, we perform a simulation experiment to illustrate the potential advantages of using our Bayesian model selection approach. We consider 5 values of p ranging from 300 to 1500, with $n = p/3$. For each fixed p , the Cholesky factor L_0 of the true concentration matrix, and the corresponding dataset, are generated by the same mechanism as in Section 6.1. Then, we perform model selection on the Cholesky factor using the four procedures outlined below.

1. *Lasso-DAG with quantile based tuning*: We implement the Lasso-DAG approach in [Shojaie and Michailidis, 2010] by choosing penalty parameters (separate for each variable i) given by $\lambda_i = 2n^{-\frac{1}{2}} Z_{\frac{0.1}{2p(i-1)}}^*$, 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.
2. *ESC Metropolis-Hastings algorithm*: We implement the Rao-Blackwellized Metropolis-Hastings algorithm for the ESC prior introduced in [Lee, Lee, and Lin, 2018] for exploring the space of the

Cholesky factor. The hyperparameters and the initial states are taken as suggested in [Lee et al., 2018]. Each MCMC chain for each row of the Cholesky factor runs for 5000 iterations with a burn-in period of 2000. All the active components in L with inclusion probability larger than 0.5 are selected. We would like to point out that since the Metropolis-Hastings algorithm needs to be executed for each row of L , the procedure could be extremely time consuming, especially in higher dimensions.

3. *DAG-Wishart log-score path search*: The hierarchical DAG-Wishart prior [Cao et al., 2019] also gives us the closed form to calculate the marginal posterior up to a constant. In particular,

$$\pi(\mathcal{D}|\mathbf{Y}) = \frac{\pi(\mathcal{D})}{\pi(\mathbf{Y})(\sqrt{2\pi})^n} \frac{z_{\mathcal{D}}(U + nS, n + \boldsymbol{\alpha}(\mathcal{D}))}{z_{\mathcal{D}}(U, \boldsymbol{\alpha}(\mathcal{D}))},$$

where $z_{\mathcal{D}}(\cdot, \cdot)$ is the normalized constant in the DAG-Wishart distribution and

$$\pi(\mathcal{D}) = \prod_{(i,j):1 \leq i < j \leq p} q^{\gamma_{ij}} (1-q)^{1-\gamma_{ij}} = \prod_{i=1}^{p-1} q^{\nu_i(\mathcal{D})} (1-q)^{p-i-\nu_i(\mathcal{D})}.$$

with $q = e^{-\eta_n n}$, where $\eta_n = d_n \left(\frac{\log p n}{n}\right)^{\frac{1/2}{1+k/2}}$. Follow the simulation procedures in previous work [Cao et al., 2019]. We set the hyperparameters as $U = I_p$ and $\alpha_i(\mathcal{D}) = \nu_i(\mathcal{D}) + 10$ for $i = 1, 2, \dots, p$ and generate candidate graphs by thresholding the modified Cholesky factor of $(S + 0.5I)^{-1}$ (S is the sample covariance matrix) on a grid from 0.1 to 0.5 by 0.0001 to get a sequence of 4000 graphs. The log posterior probabilities are computed for all candidate graphs, and the graph with the highest probability is chosen. As we discussed previously, we will see in Figure 2 that for the previous DAG-Wishart model, we always end up choosing the most sparse estimator, since the graph obtained at the thresholding value 0.5 always has the highest log posterior score. Hence, we observe that the choice $q = e^{-\eta_n n}$ though could guarantee the model selection consistency, makes the posterior stuck in very small size models and we are not able to detect the true model.

4. *Spike and slab Cholesky with beta-mixture prior/multiplicative prior*: For our Bayesian approach with spike and slab Cholesky prior and beta-mixture/multiplicative prior on the sparsity pattern of L , we adopt the similar procedure as DAG-Wishart log-score path search method. We construct two candidate sets as follows.

- (a) All the Cholesky factors with respect to the graphs on the solution paths for Lasso-DAG, CSCS and DAG-Wishart are included in our Cholesky factor candidate set.

- (b) To increase the search range, we also generate additional graphs by thresholding the modified Cholesky factor of $(S + 0.5I)^{-1}$ (S is the sample covariance matrix) on a grid from 0.1 to 0.5 by 0.0001 to get a sequence of 4000 additional Cholesky factors, and include them in the candidate set. We then search around all the above candidates using Shotgun Stochastic Search Algorithm in [Shin et al., 2018] to generate even more candidate Cholesky factors. In particular, the authors in [Shin et al., 2018] claim that the simplified algorithm can significantly lessen the simulation runtime and increase the model selection performance.

The log posterior probabilities are computed for all Cholesky factors in the candidate sets using (16), and the one with the highest probability is chosen. In Figure 2, we plot the log of marginal posterior densities under the spike and slab Cholesky prior and the multiplicative/beta-mixture prior for all the Cholesky factors under different thresholding values compared with the marginal posteriors under previous DAG-Wishart model. Unlike the DAG-Wishart distribution always favor the most sparse Cholesky factor corresponding to the largest thresholding value, we observe the maximum log posterior score occurs in the middle of the curve for our proposed models, which leads to the significant improvement of the model selection results shown in Table 1 and Table 2.

The model selection performance of these four methods is then compared using several different measures of structure such as positive predictive value, true positive rate and mathews correlation coefficient (average over 20 independent repetitions). Positive Predictive Value (PPV) represents the proportion of true non-zero entries among all the entries detected by the given procedure, True Positive Rate (TPR) measures the proportion of true non-zero entries detected by the given procedure among all the non-zero entries from the true model. PPV and TPR are defined as

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

Mathews correlation Coefficient (MCC) is commonly used to assess the performance of binary classification methods and is defined as

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{FP} + \text{TN}) \times (\text{TP} + \text{FN}) \times (\text{TN} + \text{FP}) \times (\text{TP} + \text{FN})}},$$

where TP, TN, FP and FN correspond to true positive, true negative, false positive and false negative, respectively. Note that the value of MCC ranges from -1 to 1 with larger values corresponding to better fits (-1 and 1 represent worst and best fits, respectively). Similar to MCC, one would also like the

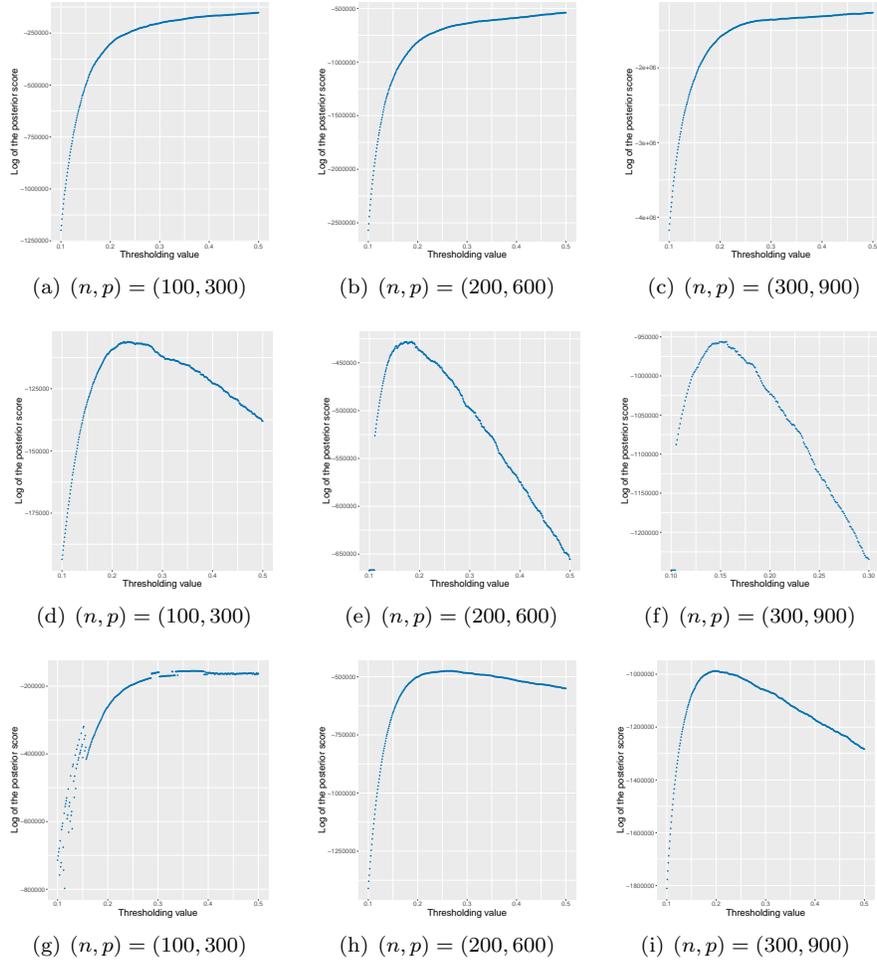


Figure 2: Log of posterior vs thresholding values under different priors. Top: DAG-Wishart; middle: Spike and slab Cholesky with beta-mixture prior; bottom: Spike and slab Cholesky with multiplicative prior.

PPV and TPR values to be as close to 1 as possible. The results are provided in Table 1 and Table 2, corresponding to different true sparsity levels. In Figure 4, we draw the heatmap comparison between the true L_0 and estimated L using our Bayesian spike and slab Cholesky approach under two different sparsity levels when $(n, p) = (100, 300)$.

p	n	Lasso-DAG			ESC			DAG-W			SSC-B			SSC-M		
		PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC
300	100	0.2	0.2	0.19	0.17	0.43	0.26	0.99	0.3	0.55	0.73	0.85	0.78	0.98	0.69	0.82
600	200	0.15	0.18	0.16	0.15	0.52	0.27	0.99	0.31	0.55	0.69	0.92	0.79	0.89	0.82	0.85
900	300	0.15	0.20	0.17	0.12	0.54	0.24	1	0.33	0.57	0.62	0.93	0.76	0.83	0.87	0.84
1200	400	0.11	0.17	0.14	0.08	0.52	0.21	1	0.33	0.58	0.61	0.94	0.76	0.78	0.90	0.84
1500	500	0.12	0.21	0.16	0.06	0.45	0.20	1	0.33	0.58	0.56	0.96	0.73	0.71	0.93	0.81

Table 1: Model selection performance table with sparsity 3%. DAG-W: DAG-Wishart log-score path search; SSC-B: Spike and slab Cholesky with beta-mixture prior; SSC-M: Spike and slab Cholesky with multiplicative prior.

p	n	Lasso-DAG			ESC			DAG-W			SSC-B			SSC-M		
		PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC	PPV	TPR	MCC
300	100	0.19	0.1	0.13	0.14	0.33	0.19	0.99	0.3	0.54	0.66	0.81	0.73	0.99	0.43	0.65
450	150	0.12	0.09	0.1	0.11	0.35	0.18	1	0.29	0.53	0.63	0.86	0.73	0.93	0.72	0.82
600	200	0.12	0.09	0.1	0.10	0.38	0.18	1	0.3	0.55	0.57	0.89	0.71	0.87	0.80	0.83
750	250	0.09	0.08	0.08	0.08	0.36	0.16	1	0.31	0.55	0.59	0.9	0.72	0.80	0.86	0.83
900	300	0.11	0.09	0.09	0.05	0.31	0.13	0.99	0.31	0.55	0.56	0.92	0.72	0.77	0.87	0.82

Table 2: Model selection performance table with sparsity 5%

It is clear that our hierarchical fully Bayesian approach with beta-mixture prior and multiplicative prior outperforms the penalized likelihood approaches, the Bayesian DAG-Wishart and ESC approach based on almost all measures. The PPV values for our Bayesian spike and slab Cholesky approach are all above 0.55, while the ones for the penalized likelihood approach and ESC are below 0.2. Though the PPV for the DAG-Wishart approach is almost 1, it is actually a consequence of the maximized log score occurring at the most sparse model. Hence, The precision (PPV) for the DAG-Wishart method is rather high, as the resulting L is extremely sparse and all the remaining non-zero entries are the true elements in L_0 . The TPR values for the proposed approaches are almost all beyond 0.70, while the ones for the penalized likelihood approaches are all below 0.27. Now again under this measure, as a result of the final sparse estimator, DAG-Wishart Bayesian approach performs very poorly compared to the spike and slab approach with beta-mixture/multiplicative prior. For the most comprehensive measure of MCC, our fully Bayesian approach outperforms all the other three methods under all the cases of (n, p) and two different sparsity levels.

It is also meaningful to compare the computational runtime between different methods. In Figure 3, we plot the run time comparison between our spike and slab Cholesky with beta-mixture prior/multiplicative prior and ESC. Since the marginal posterior is available in closed form (up to a constant) for the SSC with beta-mixture prior, we can see that the run time for SSC-B via thresholding coupled with stochastic

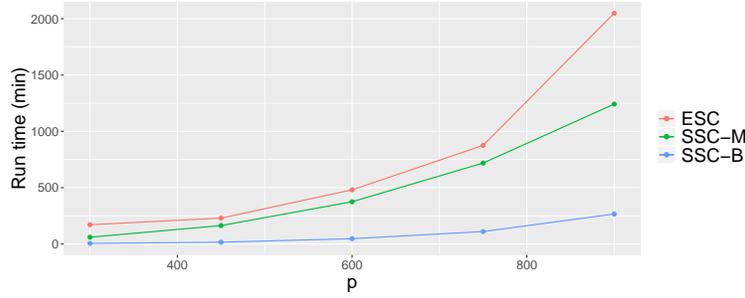


Figure 3: Run time comparison.

search is significantly lessened compared to the MCMC approach. The computational cost of ESC is extremely expensive in the sense that it requires not only additional run time, but also larger memory (more than 30GB when $p > 900$). On the other hand, for the multiplicative prior, though the model selection performance is almost the best among all the competitors, with the extra step of the Laplace approximation for calculating each posterior probability, the computational burden is quite extensive as p increases.

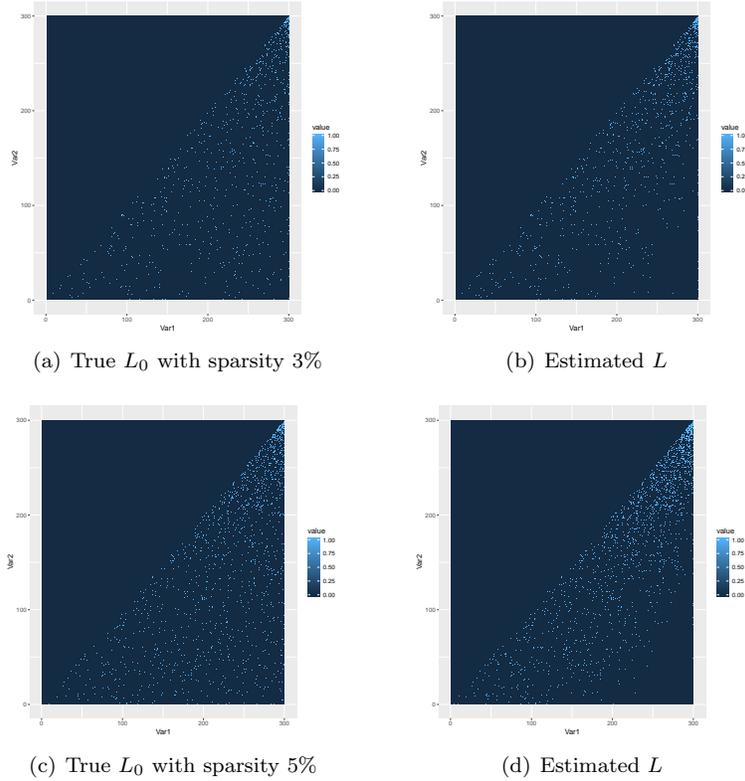


Figure 4: Heatmap comparison with $(n, p) = (100, 300)$

Overall, this experiment illustrates that the proposed hierarchical fully Bayesian approach with our

spike and slab Cholesky prior and the beta-mixture prior can be used for a broader yet computationally feasible model search, while our spike and slab Cholesky prior with the multiplicative prior though more computationally expensive, can lead to a much more significant improvement in model selection performance for estimating the sparsity pattern of the Cholesky factor and the underlying DAG.

7 Proofs

In this section, we take on the task of proving our main results presented in Theorems 4.1 to 5.3.

7.1 Proof of Theorem 4.1

The proof of Theorem 4.1 will be broken into several steps. We begin our strong selection consistency proof by first proving the Lemma 3.1 and Lemma 3.2 which give the upper bound for the prior ratio between any “non-true” model Z and the true model Z_0 .

Proof of Lemma 3.1. First note that following from model (9) and (10), we have

$$\begin{aligned}
\pi(Z) &= \int \prod_{j=1}^p \pi(\omega_j) \pi(Z|\omega_1, \dots, \omega_p) d\omega_1 \dots d\omega_p \\
&= \int \prod_{1 \leq j < k \leq p} (\omega_k \omega_j)^{Z_{kj}} (1 - \omega_k \omega_j)^{1 - Z_{kj}} \prod_{j=1}^p \pi(\omega_j) d\omega_1 \dots d\omega_p \\
&\leq \int \prod_{1 \leq j < k \leq p} (\omega_k \omega_j)^{Z_{kj}} \prod_{j=1}^p \pi(\omega_j) d\omega_1 \dots d\omega_p \\
&\leq \prod_{j=1}^p \int \omega_j^{|Z_j|} \omega_j^{\alpha_1 - 1} (1 - \omega_j)^{\alpha_2 - 1} \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \Gamma(\alpha_2)} d\omega_j \\
&\leq \prod_{j=1}^p \frac{\Gamma(\alpha_1 + \alpha_2) \Gamma(\alpha_1 + |Z_j|)}{\Gamma(\alpha_1 + \alpha_2 + |Z_j|) \Gamma(\alpha_1)}. \tag{17}
\end{aligned}$$

Denote $A_j = \left\{ \omega_j : \omega_j < \frac{\alpha_1}{\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}} \right\}$. Note that on A_j , $1 - \omega_i \omega_j > 1 - \frac{\alpha_1^2}{\max\{p^c, d^{\frac{2c}{c-2}}\}}$. Hence, by $c > 2$,

$$\begin{aligned}
\prod_{1 \leq j < k \leq p} (1 - \omega_k \omega_j)^{1 - Z_{kj}} &\geq \left(1 - \frac{\alpha_1^2}{\max\{p^c, d^{\frac{2c}{c-2}}\}} \right)^{p^2} \\
&\geq \left(1 - \frac{\alpha_1^2}{p^2} \right)^{p^2} \\
&\geq e^{-2\alpha_1^2}, \quad \text{for } p \geq \sqrt{2}\alpha_1.
\end{aligned}$$

The last inequality follows from $\frac{\log(1-x)}{x} \geq -2$, for $0 \leq x < \frac{1}{2}$. Hence, for $p \geq \sqrt{2}\alpha_1$, we have

$$\begin{aligned}
\pi(Z_0) &= \int \pi(Z_0|\omega_1, \dots, \omega_p) \prod_{j=1}^p \pi(\omega_j) d\omega_1 \dots d\omega_p \\
&\geq e^{-2\alpha_1^2} \prod_{j=1}^p \int_{A_j} \omega_j^{|Z_{0j}|+\alpha_1-1} (1-\omega_j)^{\alpha_2-1} \frac{\Gamma(\alpha_1+\alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} d\omega_j \\
&\geq e^{-2\alpha_1^2} \prod_{j=1}^p \frac{\Gamma(\alpha_1+\alpha_2)\Gamma(\alpha_1+|Z_{0j}|)}{\Gamma(\alpha_1+\alpha_2+|Z_{0j}|)\Gamma(\alpha_1)} P\left(B_j < \frac{\alpha_1}{\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}}\right), \tag{18}
\end{aligned}$$

where $B_j \sim \text{Beta}(\alpha_1 + |Z_{0j}|, \alpha_2)$. By Markov's inequality and $\alpha_2 \sim \max\{p^c, d^{\frac{2c}{c-2}}\}$, where $c > 2$, we have

$$\begin{aligned}
P\left(B_j < \frac{\alpha_1}{\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}}\right) &\geq 1 - \frac{E(B_j)}{\frac{\alpha_1}{\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}}} \\
&\geq 1 - \frac{\alpha_1 + |Z_{0j}|}{\alpha_1 (\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\})} \\
&\geq e^{-\frac{2(\alpha_1 + |Z_{0j}|)}{\alpha_1 \max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}}}, \tag{19}
\end{aligned}$$

for $p \geq 4 + \frac{4}{\alpha_1}$. The last inequality follows from

$$\begin{aligned}
\frac{\alpha_1 + |Z_{0j}|}{\alpha_1 (\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\})} &\leq \frac{\alpha_1 + d}{\alpha_1 (\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\})} \\
&\leq \frac{1}{p} + \frac{d}{\alpha_1 (\max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\})} \\
&\leq \frac{1}{p} + \frac{d}{\alpha_1 p (d^{2/(c-2)})^{c/2-1}} \\
&\leq \frac{1}{p} + \frac{1}{\alpha_1 p} \leq \frac{1}{2},
\end{aligned}$$

for $p \geq 4 + \frac{4}{\alpha_1}$. It then follows by (18) and (19) that

$$\begin{aligned}
\pi(Z_0) &\geq e^{-2\alpha_1^2} e^{-\frac{2p(\alpha_1+d)}{\alpha_1 \max\{p^{\frac{c}{2}}, d^{\frac{c}{c-2}}\}}} \prod_{j=1}^p \frac{\Gamma(\alpha_1+\alpha_2)\Gamma(\alpha_1+|Z_{0j}|)}{\Gamma(\alpha_1+\alpha_2+|Z_{0j}|)\Gamma(\alpha_1)} \\
&\geq e^{-2\alpha_1^2-2\alpha_1-\frac{2}{\alpha_2}} \prod_{j=1}^p \frac{\Gamma(\alpha_1+\alpha_2)\Gamma(\alpha_1+|Z_{0j}|)}{\Gamma(\alpha_1+\alpha_2+|Z_{0j}|)\Gamma(\alpha_1)}, \tag{20}
\end{aligned}$$

for $p \geq 4 + \frac{4}{\alpha_1} + 2\sqrt{\alpha_1}$.

Therefore, by (17) and (20) that

$$\frac{\pi(Z)}{\pi(Z_0)} \leq e^{2\alpha_1^2 + 2\alpha_1 + \frac{2}{\alpha_2}} \prod_{j=1}^p \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)}, \quad (21)$$

for $p \geq 4 + \frac{4}{\alpha_1} + 2\sqrt{\alpha_1}$. \square

Next, we prove the result on the upper bound for the marginal posterior ratio that is Lemma 3.2.

Proof of Lemma 3.2. Next, it follows model (6) to (8) that

$$\begin{aligned} & \pi(Z|\mathbf{Y}) \\ &= \int \frac{\pi(\mathbf{Y}|Z, (L, D))\pi(L|D, Z)\pi(Z)\pi(D)}{\pi(\mathbf{Y})} dLdD \\ &= \frac{\pi(Z)}{\pi(\mathbf{Y})} \int \pi(\mathbf{Y}|Z, (L, D))\pi(L|D, Z)\pi(D)dLdD. \end{aligned} \quad (22)$$

Note that

$$\begin{aligned} & \pi(\mathbf{Y}|Z, (L, D))\pi(L|D, Z)\pi(D) \\ &= \prod_{i=1}^n \left((2\pi)^{-\frac{p}{2}} \prod_{j=1}^p d_j^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{Y}_i^T (LD^{-1}L^T) \mathbf{Y}_i \right\} \right) \\ & \quad \times \prod_{j=1}^{p-1} \prod_{k=j+1}^p (N(\mathbf{0}, \tau^2 d_j) + (1 - Z_{kj})\delta_0(L_{kj})) \times \prod_{j=1}^p \pi(d_j) \\ & \propto \prod_{j=1}^{p-1} \left\{ d_j^{-\frac{n}{2}} \exp \left\{ -\frac{n \left(L_{\bar{Z},j}^{\geq} \right)^T S_{\bar{Z}}^{\geq j} L_{\bar{Z},j}^{\geq}}{2d_j} \right\} \right\} d_p^{-\frac{n}{2}} \exp \left\{ -\frac{n S_{pp}}{d_p} \right\} \\ & \quad \times \prod_{j=1}^{p-1} (d_j \tau^2)^{-\frac{|Z_j|}{2}} \exp \left\{ -\frac{\left(L_{\bar{Z},j}^{\geq} \right)^T L_{\bar{Z},j}^{\geq}}{\tau^2 d_j} \right\} \times \prod_{j=1}^p \pi(d_j). \end{aligned} \quad (23)$$

It now follows from

$$\left(L_{\bar{Z},j}^{\geq} \right)^T S_{\bar{Z}}^{\geq j} L_{\bar{Z},j}^{\geq} = \left(\mathbf{1}, \left(L_{\bar{Z},j}^{\geq} \right)^T \right) \times \begin{pmatrix} S_{jj} & \left(S_{\bar{Z},j}^{\geq} \right)^T \\ S_{\bar{Z},j}^{\geq} & S_{\bar{Z}}^{\geq j} \end{pmatrix} \times \left(\mathbf{1}, L_{\bar{Z},j}^{\geq} \right),$$

that

$$\begin{aligned}
& \exp \left\{ -\frac{n \left(L_{Z,j}^{\geq} \right)^T S_Z^{\geq j} L_{Z,j}^{\geq}}{2d_j} - \frac{\left(L_{Z,j}^{\geq} \right)^T L_{Z,j}^{\geq}}{\tau^2 d_j} \right\} \\
& = \exp \left\{ -\frac{\left(L_{Z,j}^{\geq} + \left(\tilde{S}_Z^{\geq j} \right)^{-1} \tilde{S}_{Z,j}^{\geq} \right)^T \tilde{S}_Z^{\geq j} \left(L_{Z,j}^{\geq} + \left(\tilde{S}_Z^{\geq j} \right)^{-1} \tilde{S}_{Z,j}^{\geq} \right)}{\frac{2d_j}{n}} \right\} \\
& \quad \times \exp \left\{ -\frac{\tilde{S}_{jj} - \left(\tilde{S}_{Z,j}^{\geq} \right)^T \left(\tilde{S}_Z^{\geq j} \right)^{-1} \tilde{S}_{Z,j}^{\geq}}{\frac{2d_j}{n}} + \frac{1}{2\tau^2 d_j} \right\},
\end{aligned} \tag{24}$$

where $\tilde{S} = S + \frac{1}{n\tau^2} I_p$.

It follows from Lemma 3.1, (22) and (23) that integrating out (L, D) gives us

$$\begin{aligned}
& \pi(Z|\mathbf{Y}) \\
& \propto \pi(Z) \prod_{j=1}^{p-1} \frac{1}{(n\tau^2)^{|Z_j|/2}} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{n}{2} - \lambda_1} |\tilde{S}_Z^{\geq j}|^{-\frac{1}{2}} \\
& = \pi(Z) \prod_{j=1}^{p-1} \frac{1}{(n\tau^2)^{|Z_j|/2}} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{n}{2} - \lambda_1} \left(|\tilde{S}_Z^{\geq i}| \tilde{S}_{j|Z_j} \right)^{-\frac{1}{2}},
\end{aligned} \tag{25}$$

in which $\tilde{S}_{j|Z_j} = \tilde{S}_{jj} - \left(\tilde{S}_{Z,j}^{\geq} \right)^T \left(\tilde{S}_Z^{\geq j} \right)^{-1} \tilde{S}_{Z,j}^{\geq}$.

Now note that we are interested in obtaining the posterior ratio. It immediately follows from (21) that, for $p \geq 4 + \frac{4}{\alpha_1} + 2\sqrt{\alpha_1}$, given the data Y , the posterior ratio for any Z compared to Z_0 can be simplified as

$$\begin{aligned}
& \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\
& = M_1 \prod_{j=1}^{p-1} (n\tau^2)^{-\frac{|Z_j| - |Z_{0j}|}{2}} \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \\
& \quad \times \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \left(\frac{\tilde{S}_{j|Z_{0j}}}{\tilde{S}_{j|Z_j}} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_{j|Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}}{\tilde{S}_{j|Z_j} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}} \right)^{\frac{n}{2} + \lambda_1} \\
& \triangleq M_1 \times PR'_j(Z, Z_0),
\end{aligned} \tag{26}$$

where $M_1 = e^{2\alpha_1^2 + 2\alpha_1 + \frac{2}{\alpha_1}}$, $\tilde{S} = S + \frac{1}{n\tau_{n,p}^2} I_p$ and $\tilde{S}_{j|Z_j} = \tilde{S}_{jj} - \left(\tilde{S}_{Z,j}^{\geq} \right)^T \left(\tilde{S}_Z^{\geq j} \right)^{-1} \tilde{S}_{Z,j}^{\geq}$. \square

Next, we show that in our setting, the sample and population covariance matrices are sufficiently

close with high probability. It follows by Lemma A.3 of [Bickel and Levina, 2008a] and Hanson-Wright inequality from [Rudelson and Vershynin, 2013] that there exists constants m_1, m_2 and δ depending on $\epsilon_{0,n}$ only such that for $1 \leq i, j \leq p$, we have:

$$\bar{P}(|S_{ij} - (\Sigma_0)_{ij}| \geq t) \leq m_1 \exp\{-m_2 n(t\epsilon_0)^2\}, |t| \leq \delta.$$

By the union-sum inequality, for a large enough c' such that $2 - m_2(c')^2/4 < 0$, we get that

$$\bar{P}\left(\|S - \Sigma_0\|_{\max} \geq c' \sqrt{\frac{\log p}{n}}\right) \leq mp^{2-m'c'^2/4} \rightarrow 0. \quad (27)$$

Define the event C_n as

$$C_n = \left\{ \|S - \Sigma_0\|_{\max} \geq c' \sqrt{\frac{\log p}{n}} \right\}. \quad (28)$$

We now analyze the behavior of $PR'_j(Z, Z_0)$ defined in (51) under different scenarios in a sequence of three lemmas (Lemmas 7.1 - 7.3). Recall that our goal is to find an upper bound for $PR'_j(Z, Z_0)$, such that the upper bound converges to 0 as $n \rightarrow \infty$. For all the following analyses, we will restrict ourselves to the event C_n^c .

Lemma 7.1. *If all the active elements in set Z_{j_0} are contained in the true model Z_j denoted as $Z_j \supset Z_{0j}$, then there exists N_1 (not depending on Z) such that for $n \geq N_1$ we have for some constant $\kappa > 1$, $PR'_j(Z, Z_0) \leq (2p)^{-\frac{\kappa}{\kappa}(|Z_j| - |Z_{0j}|)} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.1. We begin by simplifying the posterior ratio given in (51). Using the fact that $\sqrt{x + \frac{1}{4}} \leq \frac{\Gamma(x+1)}{\Gamma(x+\frac{1}{2})} \leq \sqrt{x + \frac{1}{2}}$ for $x > 0$ (see [Watson, 1959]), it follows from Assumption 6, $|Z_j| > |Z_{0j}|$, and $1 + x \leq e^x$, $1 - x \leq e^{-x}$, for $0 \leq x \leq 1$, that for a large enough constant M , and large enough n , we have

$$\begin{aligned} & \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \\ &= \frac{\Gamma(|Z_j| + \alpha_1)\Gamma(\alpha_1 + \alpha_2 + |Z_{0j}|)}{\Gamma(|Z_{0j}| + \alpha_1)\Gamma(\alpha_1 + \alpha_2 + |Z_j|)} \\ &\leq M \frac{(|Z_j| + \alpha_1)^{|Z_j| + \alpha_1}}{(|Z_{0j}| + \alpha_1)^{|Z_{0j}| + \alpha_1}} \frac{(\alpha_1 + \alpha_2 + |Z_{0j}|)^{\alpha_1 + \alpha_2 + |Z_{0j}|}}{(\alpha_1 + \alpha_2 + |Z_j|)^{\alpha_1 + \alpha_2 + |Z_j|}} \\ &\leq M (|Z_j| + \alpha_1)^{|Z_j| - |Z_{0j}|} \left(1 + \frac{|Z_j| - |Z_{0j}|}{|Z_{0j}| + \alpha_1}\right)^{|Z_{0j}| + \alpha_1} \\ &\quad \times (\alpha_1 + \alpha_2 + |Z_{0j}|)^{-(|Z_j| - |Z_{0j}|)} \left(1 - \frac{|Z_j| - |Z_{0j}|}{\alpha_1 + \alpha_2 + |Z_j|}\right)^{\alpha_1 + \alpha_2 + |Z_j|} \end{aligned}$$

$$\leq (c_1 p^c / |Z_j|)^{-(|Z_j| - |Z_{0j}|)}, \quad (29)$$

for some constant $c_1 > 0$.

Next, since $Z_j \supset Z_{0j}$, we can write $|\tilde{S}_Z^{\geq i}| = |\tilde{S}_{Z_0}^{\geq i}| |SC_{\tilde{S}_{Z_0}^{\geq i}}|$. Here $SC_{\tilde{S}_{Z_0}^{\geq i}}$ is the Schur complement of $\tilde{S}_{Z_0}^{\geq i}$, defined by

$$SC_{\tilde{S}_{Z_0}^{\geq i}} = D - B^T \left(\tilde{S}_{Z_0}^{\geq i} \right)^{-1} B$$

for appropriate sub matrices A and B of $\tilde{S}_Z^{\geq j}$. Since $\tilde{S}_Z^{\geq j} \geq \left(\frac{1}{n\tau_{n,p}^2} I_p \right)_Z^{\geq j}$, and $SC_{\tilde{S}_{Z_0}^{\geq i}}^{-1}$ is a principal submatrix of $\left(\tilde{S}_Z^{\geq j} \right)^{-1}$, the largest eigenvalue of $SC_{\tilde{S}_{Z_0}^{\geq i}}^{-1}$ is bounded above by $n\tau_{n,p}^2$. Therefore,

$$\left(\frac{|\tilde{S}_{Z_0}^{\geq i}|}{|\tilde{S}_Z^{\geq j}|} \right)^{\frac{1}{2}} = |SC_{\tilde{S}_{Z_0}^{\geq i}}^{-1}|^{1/2} \leq \left(\sqrt{n\tau_{n,p}^2} \right)^{|Z_j| - |Z_{0j}|}. \quad (30)$$

Denote $S_{j|Z_j} = S_{jj} - (S_{Z,j}^{\geq})^T (S_Z^{\geq j})^{-1} S_{Z,j}^{\geq}$. It immediately follows that

$$\tilde{S}_{i|Z_j} \geq S_{i|Z_j}. \quad (31)$$

Since we are restricting ourselves to the event C_n^c , it follows by (28) that

$$\|S_{Z_0}^{\geq i} - (\Sigma_0)_{Z_0}^{\geq i}\|_{(2,2)} \leq (|Z_{0j}| + 1) c' \sqrt{\frac{\log p}{n}}.$$

Therefore,

$$\begin{aligned} & \| (S_{Z_0}^{\geq i})^{-1} - ((\Sigma_0)_{Z_0}^{\geq i})^{-1} \|_{(2,2)} \\ &= \| (S_{Z_0}^{\geq i})^{-1} \|_{(2,2)} \| S_{Z_0}^{\geq i} - (\Sigma_0)_{Z_0}^{\geq i} \|_{(2,2)} \| ((\Sigma_0)_{Z_0}^{\geq i})^{-1} \|_{(2,2)} \\ &\leq (\| (S_{Z_0}^{\geq i})^{-1} - ((\Sigma_0)_{Z_0}^{\geq i})^{-1} \|_{(2,2)} + \frac{1}{\epsilon_0}) (|Z_{0j}| + 1) c' \sqrt{\frac{\log p}{n}}. \end{aligned} \quad (32)$$

Recall $d = \max_{1 \leq j \leq p-1} |Z_{0j}|$. By the assumption that $d \sqrt{\frac{\log p}{n}} \rightarrow 0$ and (32), for large enough n , we have

$$\| (S_{Z_0}^{\geq i})^{-1} - ((\Sigma_0)_{Z_0}^{\geq i})^{-1} \|_{(2,2)} \leq \frac{4c'}{\epsilon_0} d \sqrt{\frac{\log p}{n}} = o(1) \text{ and } \frac{1}{S_{i|Z_{0j}}} = \left[(S_{Z_0}^{\geq i})^{-1} \right]_{ii} \geq \frac{\epsilon_0}{2}. \quad (33)$$

¹For matrices A and B , we say $A \geq B$ if $A - B$ is positive semi-definite

Note that for any Z , $\|\tilde{S}_Z^{\geq j} - S_Z^{\geq j}\|_{\max} \leq \frac{1}{n\tau_{n,p}^2}$ gives us $\|\tilde{S}_{Z_0}^{\geq j} - S_{Z_0}^{\geq j}\|_{(2,2)} \leq (|Z_{0j}| + 1)\frac{1}{n\tau_{n,p}^2}$. Therefore,

$$\begin{aligned} & \|(\tilde{S}_{Z_0}^{\geq j})^{-1} - (S_{Z_0}^{\geq j})^{-1}\|_{(2,2)} \\ &= \|(\tilde{S}_{Z_0}^{\geq i})^{-1}\|_{(2,2)} \|\tilde{S}_{Z_0}^{\geq i} - S_{Z_0}^{\geq i}\|_{(2,2)} \| (S_{Z_0}^{\geq i})^{-1} \|_{(2,2)} \\ &\leq (\|(\tilde{S}_{Z_0}^{\geq j})^{-1} - (S_{Z_0}^{\geq j})^{-1}\|_{(2,2)} + \| (S_{Z_0}^{\geq i})^{-1} - ((\Sigma_0)_{Z_0}^{\geq i})^{-1} \|_{(2,2)} + \frac{1}{\epsilon_0})(|Z_{0j}| + 1)\frac{1}{n\tau_{n,p}^2}. \end{aligned} \quad (34)$$

Following from (33) and $\frac{d}{n\tau_{n,p}^2} \rightarrow 0$, for large enough n , (34) yields

$$\|(\tilde{S}_{Z_0}^{\geq j})^{-1} - (S_{Z_0}^{\geq j})^{-1}\|_{(2,2)} \leq \frac{8}{\epsilon_0} \frac{d}{n\tau_{n,p}^2} \text{ and } \frac{1}{\tilde{S}_{i|Z_{0j}}} = \left[(\tilde{S}_{Z_0}^{\geq i})^{-1} \right]_{ii} \geq \frac{\epsilon_0}{4}. \quad (35)$$

Hence, it follow from (35) and (33) that,

$$\left| \frac{1}{S_{i|Z_{0j}}} - \frac{1}{\tilde{S}_{i|Z_{0j}}} \right| \leq \frac{8}{\epsilon_0} \frac{d}{n\tau_{n,p}^2} \text{ and } |S_{i|Z_{0j}} - \tilde{S}_{i|Z_{0j}}| \leq c_1 \frac{d}{n\tau_{n,p}^2}, \quad (36)$$

where $c_1 = 64/\epsilon_0^3$ is a constant.

Further note that $nd_{0j}^{-1}S_{i|Z_j} \sim \chi_{n-|Z_j|}^2$ and $nd_{0j}^{-1}S_{i|Z_{0j}} \stackrel{d}{=} nd_{0j}^{-1}S_{i|Z_j} \oplus \chi_{|Z_j|-|Z_{0j}|}^2$ under the true model. Since $Z_{0j} \subset Z_j$, we get $S_{j|Z_{0j}} \geq S_{j|Z_j}$ and $\tilde{S}_{j|Z_{0j}} \geq \tilde{S}_{j|Z_j}$. It follows from Lemma 4.1 in Cao et al. [2018] that

$$P \left[\left| nd_{0j}^{-1}S_{j|Z_j} - (n - |Z_j|) \right| > \sqrt{(n - |Z_j|) \log p} \right] \leq 2p^{-\frac{1}{8}} \rightarrow 0, \quad (37)$$

and

$$P \left[\left| nd_{0j}^{-1}S_{j|Z_{0j}} - nd_{0j}^{-1}S_{j|Z_j} - (|Z_j| - |Z_{0j}|) \right| > \sqrt{(|Z_j| - |Z_{0j}|) \log p} \right] \leq 2p^{-\frac{1}{8}} \rightarrow 0. \quad (38)$$

Following from Assumption 4, Assumption 5, Lemma 3.2, (30), (31) and (35), for larger enough $n > N_1$, we have

$$\begin{aligned} & PR'_j(Z, Z_0) \\ &\leq (c_1 \sqrt{n\tau^2} p^c / |Z_j|)^{-(|Z_j| - |Z_{0j}|)} \left(1 + \frac{nd_{0j}^{-1}S_{j|Z_{0j}} - nd_{0j}^{-1}S_{j|Z_{0j}} + c_1 \frac{d}{d_{0j}\tau_{n,p}^2}}{nd_{0j}^{-1}S_{j|Z_j}} \right)^{\frac{1}{2}} \end{aligned}$$

$$\begin{aligned}
& \times \left(1 + \frac{nd_{0j}^{-1}S_{j|Z_{0j}} - nd_{0j}^{-1}S_{j|Z_j} + c_1 \frac{d}{d_{0j}\tau_{n,p}^2} + \frac{2\lambda_2}{d_{0j}}}{nd_{0j}^{-1}S_{j|Z_j} + \frac{2\lambda_2}{d_{0j}}} \right)^{\frac{n}{2} + \lambda_1} \\
& \leq (2p^{-c})^{|Z_j| - |Z_{0j}|} \left(\sqrt{\frac{\tau^2}{n}} \log n \right)^{-(|Z_j| - |Z_{0j}|)} \\
& \quad \times \exp \left\{ \frac{|Z_j| - |Z_{0j}| + \sqrt{(|Z_j| - |Z_{0j}|) \log p} + c_1 \frac{d}{\tau_{n,p}^2}}{n - |Z_j| - \sqrt{(n - |Z_j|) \log p}} \times \left(\frac{n+1}{2} + \lambda_1 \right) \right\} \quad (39) \\
& \leq (2p)^{-\frac{c}{\kappa}(|Z_j| - |Z_{0j}|)}, \text{ for some constant } \kappa > 1.
\end{aligned}$$

The second inequality follows from $\frac{d}{\tau_{n,p}^2 \log p} \rightarrow 0$, as $n \rightarrow \infty$. \square

Lemma 7.2. *If all the active elements in set Z_j are contained in the true model Z_{0j} denoted as $Z_j \subset Z_{0j}$, then there exists N_2 (not depending on Z) such that for $n \geq N_2$ we have $PR'_j(Z, Z_0) \leq p^{-\frac{2c}{\kappa}d} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.5. Now we move to discuss the scenario when Z_j is a subset of Z_{0j} , i.e., $Z_j \subset Z_{0j}$. By the similar arguments in (29), it follows from Assumption 7 and $|Z_j| < |Z_{0j}|$, that for a large enough constant c_1 and large enough n , we have

$$\begin{aligned}
& \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \\
& = \frac{\Gamma(|Z_j| + \alpha_1)\Gamma(\alpha_1 + \alpha_2 + |Z_{0j}|)}{\Gamma(|Z_{0j}| + \alpha_1)\Gamma(\alpha_1 + \alpha_2 + |Z_j|)} \\
& \leq (c_1 p^c / d)^{|Z_{0j}| - |Z_j|}. \quad (40)
\end{aligned}$$

It follows that $|\tilde{S}_{Z_0}^{\geq i}| = |\tilde{S}_{\tilde{Z}}^{\geq i}| |SC_{\tilde{S}_{\tilde{Z}}^{\geq i}}|$, where $SC_{\tilde{S}_{\tilde{Z}}^{\geq i}}$ denotes the Schur complement of $\tilde{S}_{\tilde{Z}}^{\geq i}$, defined by $SC_{\tilde{S}_{\tilde{Z}}^{\geq i}} = \tilde{A} - \tilde{B}^T (\tilde{S}_{\tilde{Z}}^{\geq i})^{-1} \tilde{B}$ for appropriate sub matrices \tilde{A} and \tilde{B} of $\tilde{S}_{Z_0}^{\geq i}$. Recall that d is the maximum number of nonzero entries among all the columns of Z_0 . It follows by (32) that if restrict to C_n^c , we have

$$\|(\tilde{S}_{Z_0}^{\geq i})^{-1} - ((\Sigma_0)_{Z_0}^{\geq i})^{-1}\|_{(2,2)} \leq \frac{4c'}{\epsilon_0} d \sqrt{\frac{\log p}{n}}$$

and

$$\|SC_{\tilde{S}_{\tilde{Z}}^{\geq i}}^{-1} - SC_{(\Sigma_0)_{\tilde{Z}}^{\geq i}}^{-1}\|_{(2,2)} \leq \frac{4c'}{\epsilon_0} d \sqrt{\frac{\log p}{n}},$$

for $n > N'_2$, in which $SC_{(\Sigma_0)_{\tilde{Z}}^{\geq i}}$ represents the Schur complement of $(\Sigma_0)_{\tilde{Z}}^{\geq i}$ given by $SC_{(\Sigma_0)_{\tilde{Z}}^{\geq i}} = \bar{A} - \bar{B}^T ((\Sigma_0)_{\tilde{Z}}^{\geq i})^{-1} \bar{B}$ for appropriate sub matrices \bar{A} and \bar{B} of $(\Sigma_0)_{Z_0}^{\geq i}$. Hence, there exists N''_2 such that, for

$n > N_2''$, we have

$$\begin{aligned} \left(\frac{|\tilde{S}_{Z_0}^{>i}|}{|\tilde{S}_{Z_j}^{>i}|} \right)^{\frac{1}{2}} &= |SC_{\tilde{S}_{Z_j}^{>i}}^{-1}|^{-\frac{1}{2}} \leq \left(\lambda_{\min} \left(SC_{(\Sigma_0)_{Z_j}^{>i}}^{-1} \right) - \frac{4c'}{\epsilon_0} d \sqrt{\frac{\log p}{n}} \right)^{-\frac{|Z_0| - |Z_j|}{2}} \\ &\leq \left(\frac{\epsilon_0}{2} \right)^{-\frac{|Z_0| - |Z_j|}{2}}. \end{aligned}$$

It follows from $Z_j \subset Z_{0j}$ that $\tilde{S}_{j|Z_{0j}} \leq \tilde{S}_{j|Z_j}$.

Let $K_1 = \frac{4c'}{\epsilon_0}$. By (51) and Proposition 5.2 in [Cao, Khare, and Ghosh, 2019], it follows that there exists N_2''' such that for $n \geq N_2'''$, we get

$$\begin{aligned} &PR'_j(Z, Z_0) \\ &\leq \left(\sqrt{\frac{2n\tau_{n,p}^2}{\epsilon_0}} c_1 p^c / d \right)^{|Z_{0j}| - |Z_j|} \left(\frac{\frac{1}{(\Sigma_0)_{j|Z_j}} + K_1 d \sqrt{\frac{\log p}{n}} - \frac{1}{n\tau_{n,p}^2}}{\frac{1}{(\Sigma_0)_{j|Z_{0j}}} - K_1 d \sqrt{\frac{\log p}{n}} - \frac{1}{n\tau_{n,p}^2}} \right)^{\frac{n}{2} + \lambda_1} \\ &\leq \left(\exp \left\{ \frac{d \log \left(\frac{2}{\epsilon_0} \right)}{n + 2\alpha_1} + \frac{2 \log \left(p^c \sqrt{n\tau_{n,p}^2} \right) (|Z_{0j}| - |Z_j|)}{n + 2\alpha_1} \right\} \right)^{\frac{n+2\lambda_1}{2}} \\ &\quad \times \left(1 + \frac{\left(\frac{1}{(\Sigma_0)_{j|Z_{0j}}} - \frac{1}{(\Sigma_0)_{j|Z_j}} \right) - 2K_1 d \sqrt{\frac{\log p}{n}}}{\frac{1}{(\Sigma_0)_{j|Z_j}} + K_1 d \sqrt{\frac{\log p}{n}}} \right)^{-\frac{n+2\lambda_1}{2}} \\ &\leq \left(\exp \left\{ \frac{d \log \left(\frac{2}{\epsilon_0} \right)}{n + 2\lambda_1} + \frac{2 \log \left(p^c \sqrt{n\tau_{n,p}^2} \right) (|Z_{0j}| - |Z_j|)}{n + 2\lambda_1} \right\} \right)^{\frac{n+2\lambda_1}{2}} \\ &\quad \times \left(1 + \frac{\epsilon_0 s_n^2 (|Z_{0j}| - |Z_j|) - 2K_1 d \sqrt{\frac{\log p}{n}}}{2/\epsilon_0} \right)^{-\frac{n+2\lambda_1}{2}}. \end{aligned} \tag{41}$$

It follows from $\frac{d \log p + d \log(n\tau_{n,p}^2)}{ns_n^2} \rightarrow 0$ and $\frac{d \sqrt{\frac{\log p}{n}}}{s_n^2} \rightarrow 0$ as $n \rightarrow \infty$, and $e^x \leq 1 + 2x$ for $x < \frac{1}{2}$, that there exists N_2'''' such that for $n \geq N_2''''$,

$$\frac{\epsilon_0 s_n^2 (|Z_{0j}| - |Z_j|) - 2K_1 d \sqrt{\frac{\log p}{n}}}{2/\epsilon_0} \geq \frac{\epsilon_0 s_n^2}{2}$$

and

$$\exp \left\{ \frac{d \log \left(\frac{2}{\epsilon_0} \right)}{n + 2\lambda_1} + \frac{2 \log \left(p^c \sqrt{n\tau_{n,p}^2} \right) (|Z_{0j}| - |Z_j|)}{n + 2\lambda_1} \right\} \leq 1 + \frac{\epsilon_0^2 s_n^2}{8}.$$

Hence, by (41), we have

$$PR'_j(Z, Z_0) \leq \left(\frac{1 + \frac{\epsilon_0^2}{8} s_n^2}{1 + \frac{\epsilon_0^2}{4} s_n^2} \right)^{\frac{n+2\lambda_1}{2}},$$

for $n \geq \max(N'_2, N''_2, N'''_2, N''''_2)$. Since there exists at least one $(L_0)_{ji}$ ($j+1 \leq i \leq p$), such that $s_n^2 \leq (L_0)_{ji}^2 \leq \frac{(\Omega_0)_{jj}}{\epsilon_0} \leq \frac{1}{\epsilon_0^2}, \epsilon_0^2 s_n^2 \leq 1$ and $e^{-x} \geq 1 - x$ when $x \geq 0$, we have for all $n \geq N_2 \triangleq \max(N'_2, N''_2, N'''_2, N''''_2)$,

$$\begin{aligned} PR'_j(Z, Z_0) &\leq \left(1 - \frac{\frac{\epsilon_0^2}{8} s_n^2}{1 + \frac{\epsilon_0^2}{4} s_n^2} \right)^{\frac{n+2\lambda_1}{2}} \leq \exp \left\{ - \left(\frac{\frac{\epsilon_0^2}{8} s_n^2}{1 + \frac{\epsilon_0^2}{4} s_n^2} \right) \left(\frac{n+2\lambda_1}{2} \right) \right\} \\ &\leq e^{-\frac{1}{10} \epsilon_0^2 s_n^2 \left(\frac{n+2\lambda_1}{2} \right)} \leq p^{-\frac{2\epsilon}{\kappa} d}, \end{aligned} \quad (42)$$

following from $\frac{d \log p}{n s_n^2} \rightarrow 0$, as $n \rightarrow \infty$. \square

Lemma 7.3. *If all the active elements in set Z_j are not contained in the true model Z_{0j} and all the active elements in set Z_{0j} are not contained in the true model Z_j , denoted as $Z_{0j} \neq Z_j$, $Z_{0j} \not\subseteq Z_j$, and $Z_{0j} \not\supseteq Z_j$, then there exists N_3 (not depending on Z) such that for $n \geq N_3$ we have $PR'_j(Z, Z_0) \leq (2p)^{-\frac{\epsilon}{\kappa} |Z_j|} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.3. Let Z^* be an arbitrary 0-1 matrix satisfying $Z_j^* = Z_j \cap Z_{0j}$. Immediately we get $pa_i(\mathcal{D}^*) \subset pa_i(\mathcal{D}_0)$ and $pa_i(\mathcal{D}^*) \subset pa_i(\mathcal{D})$. It follows from (51) that

$$\begin{aligned} PR'_j(Z, Z_0) &= (n\tau^2)^{-\frac{|Z_j| - |Z_{0j}|}{2}} \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \\ &\quad \times \left(\frac{\tilde{S}_j|_{Z_{0j}}}{\tilde{S}_j|_{Z_j}} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_j|_{Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}}{\tilde{S}_j|_{Z_j} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}} \right)^{\frac{n}{2} + \lambda_1} \\ &\leq (n\tau^2)^{-\frac{|Z_j| - |Z_j^*|}{2}} \frac{B(\alpha_1 + |Z_j|, \alpha_2)}{B(\alpha_1 + |Z_j^*|, \alpha_2)} \frac{|\tilde{S}_{Z^*}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \\ &\quad \times \left(\frac{\tilde{S}_j|_{Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}}{\tilde{S}_j|_{Z_j^*} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}} \right)^{\frac{n}{2} + \lambda_1} \\ &\quad \times (n\tau^2)^{-\frac{|Z_j^*| - |Z_{0j}|}{2}} \frac{B(\alpha_1 + |Z_j^*|, \alpha_2)}{B(\alpha_1 + |Z_{0j}|, \alpha_2)} \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_{Z^*}^{\geq j}|^{\frac{1}{2}}} \\ &\quad \times \left(\frac{\tilde{S}_j|_{Z_j^*} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}}{\tilde{S}_j|_{Z_j} - \frac{1}{n\tau_{n,p}^2} + \frac{2\lambda_2}{n}} \right)^{\frac{n}{2} + \lambda_1} \\ &\leq PR'_j(Z, Z^*) \times PR'_j(Z^*, Z_0). \end{aligned} \quad (43)$$

Note that $Z_j^* \subset Z_j$. It follows from (39) that

$$PR'_j(Z, Z^*) \leq (2p)^{-\frac{c}{\kappa}(|Z_j| - |Z_j^*|)}, \text{ for some } \kappa > 1 \text{ and } n \geq N_4. \quad (44)$$

By (57) and $Z_j^* \subset Z_{0j}$, we have

$$PR'_j(Z^*, Z_0) \leq p^{-\frac{2c}{\kappa}d}, \text{ for } n \geq N_5. \quad (45)$$

It follows from (58) and $|Z_j^*| < d$ that

$$PR'_j(Z, Z_0) \leq (2p)^{-\frac{c}{\kappa}|Z_j| - |Z_j^*|} p^{-\frac{2c}{\kappa}d} < (2p)^{-\frac{c}{\kappa}|Z_j|}, \text{ for } n \geq N_3 = \max\{N_1, N_2\}. \quad (46)$$

□

The result of Theorem 4.1 immediately follows from Lemma 7.1 to Lemma 7.3, by noting that if $Z \neq Z_0$, then there exists at least one j , such that $Z_j \neq Z_{0j}$. It follows that if we restrict to C_n^c , then

$$\max_{Z \neq Z_0} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \leq \max_{Z \neq Z_0} \prod_{j=1}^p PR_j(Z, Z_0) \rightarrow 0, \text{ as } n \rightarrow \infty, \quad (47)$$

which completes our proof of Theorem 4.1.

7.2 Proof of Theorem 4.2

We now move on to the proof of Theorem 4.2. By Lemmas 7.1 - 7.3, it follows that if we restrict to C_n^c , then for large enough constant $N > N_3$, we have

$$\begin{aligned} & \frac{1 - \pi(Z_0|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\ &= \sum_{Z \neq Z_0} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\ &\leq \sum_{j=1}^{p-1} \sum_{Z_j \neq Z_{0j}} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\ &\leq \sum_{j=1}^{p-1} \left(\sum_{Z_j \subset Z_{0j}} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} + \sum_{Z_j \supset Z_{0j}} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} + \sum_{Z_j \not\subset Z_{0j}} \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \right) \\ &\leq \sum_{j=1}^{p-1} \left(\sum_{|Z_j|=1}^{|Z_{0j}|-1} \binom{|Z_{0j}|}{|Z_j|} p^{-\frac{2c}{\kappa}d} + \sum_{|Z_j|=|Z_{0j}|}^{R_n} \binom{p - |Z_{0j}|}{|Z_j| - |Z_{0j}|} (2p)^{-\frac{c}{\kappa}(|Z_j| - |Z_j^*|)} \right) \end{aligned}$$

$$+ \sum_{|Z_j|=1}^{R_n} \binom{p}{|Z_j|} (2p)^{-\frac{c}{\kappa}|Z_j|}. \quad (48)$$

Further note that the upper bound of the binomial coefficient satisfies $\binom{p}{k} \leq p^k$, for any $1 \leq k \leq p$. It follows that when $c > 2\kappa$ for some $\kappa > 1$,

$$\frac{1 - \pi(Z_0|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Therefore, $\pi(Z_0|\mathbf{Y}) \rightarrow 1$, as $n \rightarrow \infty$, which completes our proof of the strong model selection result in Theorem 4.2.

7.3 Proof of Theorems 5.2 and 5.3

The proof of Theorem 5.2 will also be broken into several steps. We begin proving our posterior ratio consistency result by first proving the Lemma 5.1 which gives the closed form of the marginal posterior density up to a constant.

Proof of Lemma 5.1. Note that following from model (9) and (10), under the beta-mixture prior, we have

$$\begin{aligned} \pi(Z) &= \int \pi(q) \prod_{(j,k):1 \leq j < k \leq p} q^{Z_{kj}} (1-q)^{1-Z_{kj}} dq \\ &\propto \int \prod_{j=1}^{p-1} q^{\alpha_1 + |Z_j| - 1} (1-q)^{\alpha_2 + p - j - |Z_j| - 1} dq \\ &\propto B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right), \end{aligned} \quad (49)$$

where

$$\begin{aligned} &B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right) \\ &= \frac{\Gamma(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|) \Gamma(\alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j|)}{\Gamma((\alpha_1 + \alpha_2)(p-1) + \frac{p(p-1)}{2})}. \end{aligned}$$

Similar to the argument in (22) and (23), integrating out (L, D) gives us

$$\pi(Z|\mathbf{Y})$$

$$\begin{aligned}
& \propto B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right) \\
& \quad \times \prod_{j=1}^{p-1} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{n}{2}-\lambda_1} \frac{|\tilde{S}_{Z_j}^{>j}|^{-\frac{1}{2}}}{(n\tau^2)^{|Z_j|/2}} \\
& = B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right) \\
& \quad \times \prod_{j=1}^{p-1} \left(\frac{n\tilde{S}_{j|Z_j}}{2} - \frac{1}{2\tau^2} + \lambda_2 \right)^{-\frac{n}{2}-\lambda_1} \frac{\left(|\tilde{S}_{Z_j}^{>i} \tilde{S}_{j|Z_j} \right)^{-\frac{1}{2}}}{(n\tau^2)^{|Z_j|/2}}, \tag{50}
\end{aligned}$$

in which $\tilde{S}_{j|Z_j} = \tilde{S}_{jj} - (\tilde{S}_{Z_j}^{>j})^T (\tilde{S}_{Z_j}^{>j})^{-1} \tilde{S}_{Z_j}^{>j}$. \square

Now we are interested in obtaining the posterior ratio. It immediately follows from Lemma 5.1 that, given the data \mathbf{Y} , the posterior ratio for any Z compared to Z_0 can be simplified as

$$\begin{aligned}
& \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\
& = \frac{B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right)}{B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_{0j}|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_{0j}| \right)} \\
& \quad \times \prod_{j=1}^{p-1} (n\tau^2)^{-\frac{|z_j|-|z_{0j}|}{2}} \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_{Z_j}^{\geq j}|^{\frac{1}{2}}} \left(\frac{\tilde{S}_{j|Z_{0j}}}{\tilde{S}_{j|Z_j}} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_{j|Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \lambda_2}{\tilde{S}_{j|Z_j} - \frac{1}{n\tau_{n,p}^2} + \lambda_2} \right)^{\frac{n}{2}+\lambda_1}. \tag{51}
\end{aligned}$$

We begin by simplifying the posterior ratio given in (51). Using the fact that $\sqrt{x + \frac{1}{4}} \leq \frac{\Gamma(x+1)}{\Gamma(x+\frac{1}{2})} \leq \sqrt{x + \frac{1}{2}}$ for $x > 0$ (see [Watson, 1959]), it follows from Assumption 4, and $1+x \leq e^x$, $1-x \leq e^{-x}$, for $0 \leq x \leq 1$, that for a large enough constant M , and large enough n , we have

$$\begin{aligned}
& \frac{B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j| \right)}{B \left(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_{0j}|, \alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_{0j}| \right)} \\
& = \frac{\Gamma(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_j|) \Gamma(\alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_j|)}{\Gamma(\alpha_1(p-1) + \sum_{j=1}^{p-1} |Z_{0j}|) \Gamma(\alpha_2(p-1) + \frac{p(p-1)}{2} - \sum_{j=1}^{p-1} |Z_{0j}|)} \\
& \leq \prod_{j=1}^{p-1} M^{||Z_j|-|Z_{0j}||} \left(\frac{\alpha_2 + p/2 - \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}}{\alpha_1 + \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}} \right)^{-(|Z_j|-|Z_{0j}|)}, \tag{52}
\end{aligned}$$

for some constant $M > 0$.

Therefore, the posterior ratio in (51) can be bounded above by

$$\begin{aligned}
& \frac{\pi(Z|\mathbf{Y})}{\pi(Z_0|\mathbf{Y})} \\
& \leq \prod_{j=1}^{p-1} M^{\|Z_j|-|Z_{0j}\|} \left(\frac{\alpha_1 + \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}}{\alpha_2 + p/2 - \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}} \right)^{-(|Z_j|-|Z_{0j}|)} (n\tau^2)^{-\frac{|Z_j|-|Z_{0j}|}{2}} \\
& \quad \times \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \left(\frac{\tilde{S}_j|Z_{0j}}{\tilde{S}_j|Z_j} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_j|Z_{0j} - \frac{1}{n\tau_{n,p}^2} + \lambda_2}{\tilde{S}_j|Z_j - \frac{1}{n\tau_{n,p}^2} + \lambda_2} \right)^{\frac{n}{2} + \lambda_1} \\
& \triangleq PR_j(Z, Z_0). \tag{53}
\end{aligned}$$

We now analyze the behavior of $PR_j(Z, Z_0)$ defined in (51) under different scenarios in a sequence of three lemmas (Lemmas 7.4 - 7.6). Recall that our goal is to find an upper bound for $PR_j(Z, Z_0)$, such that the upper bound converges to 0 as $n \rightarrow \infty$. For all the following analyses, we will restrict ourselves to the event C_n^c .

Lemma 7.4. *If all the active elements in set Z_{j_0} are contained in the true model Z_j denoted as $Z_j \supset Z_{0j}$, then there exists N_4 (not depending on Z) such that for $n \geq N_4$ we have for some constant $\kappa > 1$, $PR_j(Z, Z_0) \leq (2p)^{-\frac{\max\{c,1\}}{\kappa}(|Z_j|-|Z_{0j}|)} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.4. We begin by simplifying the posterior ratio given in (51). It follows from Assumption 7, $|Z_j| > |Z_{0j}|$, that for a large enough constant M , and large enough n , we have

$$\begin{aligned}
& M^{\|Z_j|-|Z_{0j}\|} \left(\frac{\alpha_2 + p/2 - \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}}{\alpha_1 + \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}} \right)^{-(|Z_j|-|Z_{0j}|)} \\
& \leq (c_1 p^{\max\{c,1\}}/n)^{-(|Z_j|-|Z_{0j}|)}, \tag{54}
\end{aligned}$$

for some constant $c_1 > 0$.

Following the similar arguments leading up to (55), by Assumption 4, Assumption 5, for larger enough $n \geq N_4$, we have

$$\begin{aligned}
& PR_j(Z, Z_0) \\
& \leq (c_1 \sqrt{n\tau^2} p^{\max\{c,1\}}/n)^{-(|Z_j|-|Z_{0j}|)} \left(1 + \frac{nd_{0j}^{-1} S_j|Z_{0j} - nd_{0j}^{-1} S_j|Z_{0j} + c_1 \frac{d}{d_{0j} \tau_{n,p}^2}}{nd_{0j}^{-1} S_j|Z_j} \right)^{\frac{1}{2}}
\end{aligned}$$

$$\begin{aligned}
& \times \left(1 + \frac{nd_0^{-1}S_{j|Z_{0j}} - nd_0^{-1}S_{j|Z_j} + c_1 \frac{d}{d_0 \tau_{n,p}^2}}{nd_0^{-1}S_{j|Z_j}} \right)^{\frac{n}{2} + \lambda_1} \\
& \leq (2p^{-\max\{c,1\}})^{|Z_j| - |Z_{0j}|} \left(\sqrt{\frac{\tau^2}{n}} \log n \right)^{-(|Z_j| - |Z_{0j}|)} \\
& \quad \times \exp \left\{ \frac{|Z_j| - |Z_{0j}| + \sqrt{(|Z_j| - |Z_{0j}|) \log p} + c_1 \frac{d}{\tau_{n,p}^2}}{n - |Z_j| - \sqrt{(n - |Z_j|) \log p}} \times \left(\frac{n+1}{2} + \lambda_1 \right) \right\} \\
& \leq (2p)^{-\frac{\max\{c,1\}}{\kappa} (|Z_j| - |Z_{0j}|)}, \text{ for some constant } \kappa > 1.
\end{aligned} \tag{55}$$

The second inequality follows from $\frac{d}{\tau_{n,p}^2 \log p} \rightarrow 0$, as $n \rightarrow \infty$. \square

Lemma 7.5. *If all the active elements in set Z_j are contained in the true model Z_{0j} denoted as $Z_j \subset Z_{0j}$, then there exists N_5 (not depending on Z) such that for $n \geq N_5$, we have $PR_j(Z, Z_0) \leq p^{-\frac{2c}{\kappa}d} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.5. Now we move to discuss the scenario when Z_j is a subset of Z_{0j} , i.e., $Z_j \subset Z_{0j}$. By the similar arguments in (52), it follows from Assumption 7 and $|Z_j| < |Z_{0j}|$, that for a large enough constant c_1 and large enough n , we have

$$\begin{aligned}
& M^{|Z_j| - |Z_{0j}|} \left(\frac{\alpha_2 + p/2 - \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}}{\alpha_1 + \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}} \right)^{-(|Z_j| - |Z_{0j}|)} \\
& \leq (c_1 p^{\max\{c,1\}})^{|Z_{0j}| - |Z_j|},
\end{aligned} \tag{56}$$

It follows from the similar arguments leading up to (42) that there exists $N_5 > 0$, such that for $n \geq N_5$,

$$PR_j(Z, Z_0) \leq p^{-\frac{2c}{\kappa}d}. \tag{57}$$

\square

Lemma 7.6. *If all the active elements in set Z_j are not contained in the true model Z_{0j} and all the active elements in set Z_{0j} are not contained in the true model Z_j , denoted as $Z_{0j} \neq Z_j$, $Z_{0j} \not\subseteq Z_j$, and $Z_{0j} \not\supseteq Z_j$, then there exists N_6 (not depending on Z) such that for $n \geq N_6$ we have $PR_j(Z, Z_0) \leq (2p)^{-\frac{\max\{c,1\}}{\kappa} (|Z_j| - |Z_j^*|) - \frac{2c}{\kappa}d} \rightarrow 0$, as $n \rightarrow \infty$.*

Proof of Lemma 7.6. Let Z^* be an arbitrary 0-1 matrix satisfying $Z_j^* = Z_j \cap Z_{0j}$. Immediately we get

$pa_i(\mathcal{D}^*) \subset pa_i(\mathcal{D}_0)$ and $pa_i(\mathcal{D}^*) \subset pa_i(\mathcal{D})$. It follows from (51) that

$$\begin{aligned}
PR_j(Z, Z_0) &\leq M^{\|Z_j| - |Z_{0j}|} \left(\frac{\alpha_1 + \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}}{\alpha_2 + p/2 - \frac{\sum_{j=1}^{p-1} |Z_j|}{p-1}} \right)^{-(|Z_j| - |Z_{0j}|)} (n\tau^2)^{-\frac{|Z_j| - |Z_{0j}|}{2}} \\
&\quad \times \frac{|\tilde{S}_{Z_0}^{\geq j}|^{\frac{1}{2}}}{|\tilde{S}_Z^{\geq j}|^{\frac{1}{2}}} \left(\frac{\tilde{S}_j|_{Z_{0j}}}{\tilde{S}_j|_{Z_j}} \right)^{\frac{1}{2}} \left(\frac{\tilde{S}_j|_{Z_{0j}} - \frac{1}{n\tau_{n,p}^2} + \lambda_2}{\tilde{S}_j|_{Z_j} - \frac{1}{n\tau_{n,p}^2} + \lambda_2} \right)^{\frac{n}{2} + \lambda_1} \\
&\leq PR_j(Z, Z^*) \times PR_j(Z^*, Z_0).
\end{aligned} \tag{58}$$

Note that $Z_j^* \subset Z_j$. It follows from (55) that

$$PR_j(Z, Z^*) \leq (2p)^{-\frac{\max\{c,1\}}{\kappa}(|Z_j| - |Z_j^*|)}, \text{ for some } \kappa > 1 \text{ and } n \geq N_4. \tag{59}$$

By (42) and $Z_j^* \subset Z_{0j}$, we have

$$PR_j(Z^*, Z_0) \leq p^{-\frac{2c}{\kappa}d}, \text{ for } n \geq N_5. \tag{60}$$

It follows from (58) and $|Z_j^*| < d$ that

$$PR_j(Z, Z_0) \leq (2p)^{-\frac{\max\{c,1\}}{\kappa}(|Z_j| - |Z_j^*|) - \frac{2c}{\kappa}d}, \text{ for } n > N_6 = \max\{N_4, N_5\}. \tag{61}$$

□

For any $Z \neq Z_0$, it follows that there exists at least one $1 \leq i \leq p-1$, such that $Z_j \neq Z_{0j}$. Hence, by (55), (42) and (61), we have

$$PR_j(Z, Z_0) \rightarrow 0, \text{ as } n \rightarrow \infty. \tag{62}$$

The results of Theorem 5.2 and 5.3 can be immediately obtained from Lemma 7.4 to Lemma 7.6 by following the same arguments leading up to (48).

8 Discussion

In this paper, we investigate the theoretical consistency properties for the high-dimensional sparse DAG models based on the spike and slab prior introduced on the Cholesky parameter and appropriate mul-

multiplicative and beta-mixture priors on the indicator probabilities. We establish both posterior ratio consistency and the strong model selection consistency under more general conditions than those in the existing literature. In particular, our consistency result requires much more relaxed conditions on the dimensionality and sparsity. In addition, rather than treating q as a constant and controlling its rate, by either putting an extra layer prior on q or placing the multiplicative prior over the space of Z , we avoid the potential issues of the model being stuck in rather sparse space. Finally, the simulation study shows that not only the proposed models yield desired asymptotic consistency, in the same time can also give a better model selection performance.

References

- B. Aragam, A. Amini, and Q. Zhou. Learning directed acyclic graphs with penalized neighbourhood regression. <https://arxiv.org/abs/1511.08963>, 2015.
- S. Banerjee and S. Ghosal. Posterior convergence rates for estimating large precision matrices using graphical models. *Electronic Journal of Statistics*, 8:2111–2137, 2014.
- S. Banerjee and S. Ghosal. Bayesian structure learning in graphical models. *Journal of Multivariate Analysis*, 136:147–162, 2015.
- E. Ben-David, T. Li, H. Massam, and B. Rajaratnam. High dimensional bayesian inference for gaussian directed acyclic graph models. *Technical Report*, <http://arxiv.org/abs/1109.4371>, 2016.
- P. J. Bickel and E. Levina. Regularized estimation of large covariance matrices. *Ann. Statist.*, 36:199–227, 2008a.
- Peter J. Bickel and Elizaveta Levina. Covariance regularization by thresholding. *Ann. Statist.*, 36(6): 2577–2604, 12 2008b. doi: 10.1214/08-AOS600. URL <https://doi.org/10.1214/08-AOS600>.
- X. Cao, K. Khare, and M. Ghosh. High-dimensional posterior consistency for hierarchical non-local priors in regression. <https://arxiv.org/abs/1709.06607>, 2018.
- X. Cao, K. Khare, and M. Ghosh. Posterior graph selection and estimation consistency for high-dimensional bayesian dag models. *Ann. Statist.*, 47(1):319–348, 02 2019.
- C. M. Carvalho and J. G. Scott. Objective Bayesian model selection in Gaussian graphical models. *Biometrika*, 96(3):497–512, 05 2009.

- N. El Karoui. Spectrum estimation for large dimensional covariance matrices using random matrix theory. *Annals of Statistics*, 36:2757–2790, 2008.
- Noureddine El Karoui. Tracy–widom limit for the largest eigenvalue of a large class of complex sample covariance matrices. *Ann. Probab.*, 35(2):663–714, 03 2007. doi: 10.1214/009117906000000917.
- J. Huang, N. Liu, M. Pourahmadi, and L. Liu. Covariance selection and estimation via penalised normal likelihood. *Biometrika*, 93:85–98, 2006.
- V. Johnson and D. Rossell. Bayesian model selection in high-dimensional settings. *J. Amer. Statist. Assoc.*, 107(498):649–660, 2012.
- K. Khare, S. Oh, S. Rahman, and B. Rajaratnam. A convex framework for high-dimensional sparse cholesky based covariance estimation in gaussian dag models. *Preprint, Department of Statistics, University of Florida*, 2017.
- K. Lee and J. Lee. Estimating large precision matrices via modified cholesky decomposition. <https://arxiv.org/abs/1707.01143>, 2017.
- Kyoungjae Lee, Jaeyong Lee, and Lizhen Lin. Minimax posterior convergence rates and model selection consistency in high-dimensional dag models based on sparse cholesky factors. *Ann. Statist., to appear*, 2018.
- N. Narisetty and X. He. Bayesian variable selection with shrinking and diffusing priors. *Ann. Statist.*, 42:789–817, 2014.
- M. Pourahmadi. Cholesky decompositions and estimation of a covariance matrix: Orthogonality of variance–correlation parameters. *Biometrika*, 94:1006–1013, 2007.
- A. J. Rothman, E. Levina, and J. Zhu. A new approach to cholesky-based covariance regularization in high dimensions. *Biometrika*, 97:539–550, 2010.
- Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18:9 pp., 2013.
- P. Rutimann and P. Bühlmann. High dimensional sparse covariance estimation via directed acyclic graphs. *Electronic Journal of Statistics*, 3:1133–1160, 2009.
- M. Shin, A. Bhattacharya, and V. Johnson. Scalable bayesian variable selection using nonlocal prior densities in ultrahigh-dimensional settings. *Statist. Sinica*, 28:1053–1078, 2018.

- A. Shojaie and G. Michailidis. Penalized likelihood methods for estimation of sparse high-dimensional directed acyclic graphs. *Biometrika*, 97:519–538, 2010.
- Linda S. L. Tan, Ajay Jasra, Maria De Iorio, and Timothy M. D. Ebbels. Bayesian inference for multiple gaussian graphical models with application to metabolic association networks. *Ann. Appl. Stat.*, 11(4): 2222–2251, 12 2017.
- G.N. Watson. A note on gamma functions. *Proc. Edinburgh Math. Soc.*, 11:7–9, 1959.
- R. Xiang, K. Khare, and M. Ghosh. High dimensional posterior convergence rates for decomposable graphical models. *Electronic Journal of Statistics*, 9:2828–2854, 2015.
- Xiaofan Xu and Malay Ghosh. Bayesian variable selection and estimation for group lasso. *Bayesian Anal.*, 10(4):909–936, 12 2015. doi: 10.1214/14-BA929.
- Yun Yang, Martin J. Wainwright, and Michael I. Jordan. On the computational complexity of high-dimensional bayesian variable selection. *Ann. Statist.*, 44(6):2497–2532, 12 2016. doi: 10.1214/15-AOS1417. URL <https://doi.org/10.1214/15-AOS1417>.
- G. Yu and J. Bien. Learning local dependence in ordered data. *arXiv:1604.07451*, 2016.