



Statistical Theory and Related Fields

ISSN: (Print) (Online) Journal homepage: https://www.tandfonline.com/loi/tstf20

Posterior propriety of an objective prior for generalized hierarchical normal linear models

Cong Lin, Dongchu Sun & Chengyuan Song

To cite this article: Cong Lin, Dongchu Sun & Chengyuan Song (2022) Posterior propriety of an objective prior for generalized hierarchical normal linear models, Statistical Theory and Related Fields, 6:4, 309-326, DOI: <u>10.1080/24754269.2021.1978206</u>

To link to this article: <u>https://doi.org/10.1080/24754269.2021.1978206</u>

© 2022 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group



0

Published online: 30 Jul 2022.

Submit your article to this journal \square

Article views: 159



💽 View related articles 🗹

🕨 View Crossmark data 🗹



OPEN ACCESS Check for updates

Posterior propriety of an objective prior for generalized hierarchical normal linear models

Cong Lin^a, Dongchu Sun^{a,b} and Chengyuan Song^c

^a School of Statistics, East China Normal University, Shanghai, People's Republic of China; ^bDepartment of Statistics, University of Nebraska-Lincoln, Lincoln, NE, USA; ^cBoehringer Ingelheim (China), Shanghai, People's Republic of China

ABSTRACT

Bayesian Hierarchical models has been widely used in modern statistical application. To deal with the data having complex structures, we propose a generalized hierarchical normal linear (GHNL) model which accommodates arbitrarily many levels, usual design matrices and 'vanilla' covariance matrices. Objective hyperpriors can be employed for the GHNL model to express ignorance or match frequentist properties, yet the common objective Bayesian approaches are infeasible or fraught with danger in hierarchical modelling. To tackle this issue, [Berger, J., Sun, D., & Song, C. (2020b). An objective prior for hyperparameters in normal hierarchical models. Journal of Multivariate Analysis, 178, 104606. https://doi.org/10.1016/j.jmva.2020.104606] proposed a particular objective prior and investigated its properties comprehensively. Posterior propriety is important for the choice of priors to guarantee the convergence of MCMC samplers. James Berger conjectured that the resulting posterior is proper for a hierarchical normal model with arbitrarily many levels, a rigorous proof of which was not given, however. In this paper, we complete this story and provide an user-friendly guidance. One main contribution of this paper is to propose a new technique for deriving an elaborate upper bound on the integrated likelihood, but also one unified approach to checking the posterior propriety for linear models. An efficient Gibbs sampling method is also introduced and outperforms other sampling approaches considerably.

1. Introduction

Bayesian hierarchical models (or multilevel models) have been extensively used in the modern application, including education (Raudenbush & Bryk, 1986), psychology (Lindenberger & Pötter, 1998), clinical trials (Xia et al., 2011), economics (Shimotsu, 2010) and many other applied statistical fields. The fundamental idea of hierarchical modelling is to think of the lowest-level units (smallest and most numerous) as organized into a hierarchy of successively higher-level units. For example, students are in classes, classes are in schools, schools are in districts, and districts are in states. Accordingly, hierarchical models are naturally applicable to the survey, observational or experimental data involved with complicated nesting. However, the most commonly used and fully discussed hierarchical models are merely of two levels. Goldstein (2011) and Berger et al. (2020b) have ever defined 3-level hierarchical model and implemented statistical analysis on that. The hierarchical model with more levels was usually avoided by the researchers for the reason of analytical difficulty and intractable computation. To the best of authors' knowledge, a general hierarchical linear model with arbitrary levels seems to have never been defined or studied. In this paper, we will introduce the definition of a *generalized hierarchical normal linear* (GHNL) model and carry out an in-depth theoretical investigation of Bayesian inference for the GHNL models.

In order to implement fully Bayesian analysis, priors are supposed to be specified on the hyperparameters (parameters at higher levels of the hierarchical model). Improper (objective) priors are often used to express ignorance or to match frequentist properties (see the review article, Consonni et al., 2018). When using improper priors, an important issue whether the resulting posterior distributions are proper arises. As Hobert and Casella (1996) stated, without proper precaution, misuse of improper priors, sometimes unknowingly, will result in practical difficulties, such as the nonconvergence of the Gibbs sampler. The enormous practical importance of posterior propriety motivates us to explore it in the framework of GHNL modelling afterwards. There is also a vast modern literature investigating the posterior propriety of improper priors applied to a large variety of models, such as, Sun et al. (2001), Speckman and Sun (2003), Berger et al. (2005) and Michalak and Morris (2016).

A great deal of efforts have been devoted to the development of objective hyperpriors in hierarchical

ARTICLE HISTORY

Received 23 January 2021 Revised 31 August 2021 Accepted 31 August 2021

Taylor & Francis

Taylor & Francis Group

KEYWORDS

Hierarchical linear model; linear mixed-effect model; objective Bayesian analysis; posterior propriety; Gibbs sampling

CONTACT Cong Lin 🖾 conglin_stat@163.com

^{© 2022} The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

modelling, such as, Daniels and Kass (1999), Everson and Morris (2000), Gelman (2006), Gustafson et al. (2006), Berger et al. (2005) and Berger et al. (2020b). Formal objective Bayesian approaches, like the Jeffreys-rule prior or reference prior are only feasible for the simple hierarchical settings. For instance, the exact Jeffreys-rule prior for covariance matrices at higher level depends on the parameters from the lower level of the model, leading to plenty of difficulties in formulation and computation. Therefore, a common way is to use less formal approaches, such as applying formal objective priors from non-hierarchical models to hierarchical modelling. Unfortunately, the non-hierarchical Jefferys-rule prior and reference prior typically yield improper posteriors in the hierarchical settings (cf. Berger et al., 2005). Those who can recognize this problem often use constant priors instead for higher level variance components. However, the constant prior is so diffuse that it requires twice as many observations as logically needed to achieve posterior propriety (cf. Berger et al., 2005 and Berger et al., 2020b). In other words, the extra observations required are wasted on correcting the over-diffuse tail of the constant prior. The most powerful tool known for detecting over-diffuse hyperpriors is by looking at the frequentist notion of admissibility of resulting estimators (see Berger et al., 2005 for discussions and references). Sensible choices of objective hyperpriors are on the boundary of admissibility, being as diffuse as possible without leading to inadmissible estimators.

Berger et al. (2005) studied the propriety and admissibility of a number of hyperprios, but no overall conclusion was reached as to a specific prior to recommend. The reasons are as follows: (a) the admissibility of the leading candidate prior was unable to get proved; (b) the proposed computation methods were only efficient for relatively low-dimensional covariance matrices and remained quite challenging for the candidate priors; (c) the hierarchical model discussed was merely of two levels, and the results are not adaptive to a general hierarchical model with many levels. To address this issue, Berger et al. (2020b) recommended a particular objective prior for use in all normal hierarchical models. Consider the following canonical form of 2-level hierarchical normal model. Suppose that, independently, $\mathbf{y}_i \sim N_k(\boldsymbol{\theta}_i, \mathbf{I}_k)$ and $\boldsymbol{\theta}_i \sim N_k(\boldsymbol{\beta}, \mathbf{V})$ for $i = 1, \dots, m$, where $N_k(\cdot, \cdot)$ denotes the *k*-dimensional normal distribution, y_i are $k \times 1$ observation vectors, θ_i are the $k \times 1$ unobserved mean vectors, $\boldsymbol{\beta}$ is a $p \times 1$ 'hypermean' vector, and $V \in \mathbb{R}^{k \times k}$ is an unknown 'hypercovariance' matrix. Berger et al. (2020b) proposed a particular combination of independent priors on hyperparameters β and V as

$$\pi(\boldsymbol{\beta}) \propto \frac{1}{(1+\|\boldsymbol{\beta}\|^2)^{(k-1)/2}},$$

$$\pi(V) \propto \frac{1}{|V|^{1-1/(2k)} \prod_{1 \le s < t \le k} (v_s - v_t)}, \quad (1)$$

where $v_1 > v_2 > \cdots > v_k > 0$ are the ordered eigenvalues of *V*. The recommendation (1) for hyperpriors was justified by Berger et al. (2020b) from the aspects of admissibility, ease of computation and performance. Most importantly, prior (1) is adapted to being used at any level in hierarchical modelling, which is not true for other proposed objective priors as previously mentioned.

Since it is hazardous to skip demonstrating propriety at the risk of making inference from improper posterior, Berger et al. (2020b) has shown the posterior propriety of a 3-level hierarchical model using prior (1), while assuming square design matrices for a technical reason. Berger et al. (2020b) also conjectured that the posterior is proper with the recommended prior being utilized at all levels of a hierarchical normal model with arbitrarily many levels, a rigorous proof of which was unable to be provided, however. In this paper, we complete this story and prove the posterior propriety for the GHNL models in general situations. Besides, as pointed out in Michalak and Morris (2016), researchers have been finding it daunting and time-consuming to inspect posterior propriety when using improper priors, except in the simplest models. For this reason, we supply a userfriendly guidance for checking posterior propriety to practitioners in different practical situations.

In Section 2, we give an explicit definition of the GHNL model which accommodates arbitrarily many levels and usual design matrices. It is important to note that we are considering the 'vanilla' covariance matrix problem herein. We are not assuming any special structures or sparsity for hypercovariance matrices. The association between the GHNL model and a linear mixed-effect model is also discussed. In Section 3, we demonstrate that recommended prior yields a proper posterior in the framework of GHNL modelling. In addition, we exhibit a guidance for checking posterior propriety. An efficient MCMC algorithm for sampling from the posterior is introduced in Section 4. Section 5 provides some concluding remarks and further generalizations.

2. Generalized hierarchical normal linear model

In this section, we will introduce the definition of a GHNL model with (r + 1) levels, where $r \ge 1$. The association between the GHNL model and a linear mixed-effect model will be demonstrated as well, which brings an insight into the GHNL model. At last, the recommended prior on the hyperparameters of the GHNL model will be presented and discussed. Firstly, we introduce some notations to be used in the main body of this paper.

Notations Let $[k] = \{1, 2, ..., k\}$ for a positive integer k; $1_{\{\cdot\}}$ stands for the indicator function; $N_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

represents the *k*-dimensional normal distribution with mean μ and covariance Σ ; $\mathcal{N}_k(\mu, \Sigma)$ denotes a *k*-dimensional normal random variable with mean μ and covariance Σ ; for a symmetric matrix A, A > (<)0 means that A is a positive (negative) definite matrix, and $A \ge (\leq)0$ denotes that A is a non-negative (non-positive) definite matrix.

2.1. Model structure

Berger et al. (2020b) proposed a 3-level hierarchical model with the form:

$$\begin{cases} \text{Level 1}: \mathbf{y}_i = \boldsymbol{\theta}_i + \mathcal{N}_k(\mathbf{0}, \mathbf{I}_k), & i \in [m];\\ \text{Level 2}: \boldsymbol{\theta}_i = \mathbf{Z}_i \boldsymbol{\beta} + \mathcal{N}_k(\mathbf{0}, \mathbf{V}), & \boldsymbol{\beta}^\top = (\boldsymbol{\beta}_1^\top, \\ \dots, \boldsymbol{\beta}_s^\top);\\ \text{Level 3}: \boldsymbol{\beta}_j = \boldsymbol{\eta} + \mathcal{N}_p(\mathbf{0}, \mathbf{W}), & j \in [s], \end{cases}$$
(2)

where y_i are $k \times 1$ observation vectors, θ_i are the $k \times 1$ unobserved mean vectors, η is a $p \times 1$ 'hypermean' vector, $V \in \mathbb{R}^{k \times k}$ and $W \in \mathbb{R}^{p \times p}$ are unknown 'hypercovariance' matrices, and Z_i are $k \times sp$ known matrices. At last, all the normal random variables in model (2) are mutually independent. Based on the 3-level hierarchical normal model, a more general hierarchical model with (r + 1) levels $(r \ge 1)$ can be constructed as

$$\begin{cases} \text{Level 1}: & y_{i_0} = Z_{0i_0} \theta_1 & i_0 \in [m_0], \\ + \mathcal{N}_{k_0}(\mathbf{0}, I_{k_0}), & \theta_1^\top = (\theta_{11}^\top, \\ \dots, \theta_{1m_1}^\top); & \dots, \theta_{1m_1}^\top); \\ \text{Level 2}: & \theta_{1i_1} = Z_{1i_1} \theta_2 & i_1 \in [m_1], \theta_2^\top \\ + \mathcal{N}_{k_1}(\mathbf{0}, \mathbf{V}_1), & \dots, \theta_{2m_2}^\top); \\ \vdots & \vdots & \vdots \\ \theta_{r-1, i_{r-1}} & i_{r-1} \in [m_{r-1}], \\ \text{Level } r: & = Z_{r-1, i_{r-1}} \theta_r & \theta_r^\top = (\theta_{r1}^\top, \\ + \mathcal{N}_{k_{r-1}}(\mathbf{0}, \mathbf{V}_{r-1}), & \dots, \theta_{rm_r}^\top); \\ \text{Level } r+1: & \theta_{ri_r} = Z_{ri_r} \eta \\ + \mathcal{N}_{k_r}(\mathbf{0}, \mathbf{V}_r), & i_r \in [m_r]. \end{cases} \end{cases}$$
(3)

Firstly, all the normal random variables noted in the above model are mutually independent. Within model (3), the output of level (j + 1) consists of m_j units whose values are $k_j \times 1$ vectors for j = 0, 1, ..., r. By stacking the output units of level (j + 1) on top of one another, we can obtain the *outcome vector* of level (j + 1) as θ_j for $j \in [r]$ and $\mathbf{y} = (\mathbf{y}_1^\top, ..., \mathbf{y}_{m_0}^\top)^\top$ for level 1. Then θ_j 's are $(m_j k_j) \times 1$ vectors and \mathbf{y} is a $(m_0 k_0) \times 1$ vector. In fact, only the outcome of the lowest level can be observed, and the outcomes of higher levels are inaccessible and latent variables. Hence, the outcome variables of interest are always situated at the lowest level of the hierarchy. Different units in the same level share in common *input effects* (intercept can be included) which

Table 1. Summary of certain important notations within model (3) and j = 0, 1, ..., r.

Notation	Meaning
r+1	Total number of levels
m _i	Number of the output units in level $(j + 1)$
k _i	Dimension of each output unit in level $(j + 1)$
d	Dimension of the fixed effect η

are exactly the outcome vectors from the upper level, except that the input effect of level (r + 1) is η which is a $d \times 1$ vector of fixed effects. In addition, the units from the same level have the same variance component. The variance component within level (j + 1) is denoted by $V_j \in \mathbb{R}^{k_j \times k_j}$ for $j \in [r]$ and accounts for the magnitude of random variation within the corresponding level. The covariance matrices V_j are unobserved for $j \in [r]$. The matrices Z_{ji_j} are $k_j \times (m_{j+1}k_{j+1})$ matrices and denote the matrices of observed covariates for unit i_j in level j, where j = 0, 1, ..., r and $i_j \in [m_j]$. It is natural to assume that there exist at least two units in each level and the dimensions of all units and η are no less than 1, mathematically, $m_i \ge 2$ and $k_i \ge 1$ for $j = 0, 1, \ldots, r$, and $d \ge 1$. Table 1 summarizes several important notations that will mainly affect the results for the posterior propriety in Section 3.

The extensions from Berger et al. (2020b)'s model (2) to model (3) are two-fold, and model (3) accommodates arbitrarily many levels and usual covariate matrices. Further define $\mathbf{Z}_j = \left\{ \mathbf{Z}_{j1}^\top, \ldots, \mathbf{Z}_{jm_j}^\top \right\}^\top$ for $j = 0, 1, \ldots, r$. Then \mathbf{Z}_j are $(m_j k_j) \times (m_{j+1} k_{j+1})$ matrices for $j = 0, 1, \ldots, (r-1), \mathbf{Z}_r$ is an $(m_r k_r) \times d$ matrix, and an alternative representation of the (r + 1)-level hierarchical model (3) is thereby given by

$$\begin{cases} \text{Level } 1: & (\boldsymbol{y}|\boldsymbol{\theta}_{1}) & \sim & N_{m_{0}k_{0}}(\boldsymbol{Z}_{0}\boldsymbol{\theta}_{1},\boldsymbol{I}_{m_{0}k_{0}}); \\ \text{Level } 2: & (\boldsymbol{\theta}_{1}|\boldsymbol{\theta}_{2},\boldsymbol{V}_{1}) & \sim & N_{m_{1}k_{1}}(\boldsymbol{Z}_{1}\boldsymbol{\theta}_{2}, \\ \boldsymbol{I}_{m_{1}}\otimes\boldsymbol{V}_{1}); \\ \vdots & \vdots & \vdots \\ \text{Level } r: & (\boldsymbol{\theta}_{r-1}|\boldsymbol{\theta}_{r},\boldsymbol{V}_{r-1}) & \sim & \frac{N_{m_{r-1}k_{r-1}}(\boldsymbol{Z}_{r-1}\boldsymbol{\theta}_{r},\boldsymbol{I}_{m_{r-1}})}{\otimes \boldsymbol{V}_{r-1};} \\ \text{Level } r+1: & (\boldsymbol{\theta}_{r}|\boldsymbol{\eta},\boldsymbol{V}_{r}) & \sim & N_{m_{r}k_{r}}(\boldsymbol{Z}_{r}\boldsymbol{\eta}, \\ \boldsymbol{I}_{m_{r}}\otimes\boldsymbol{V}_{r}). \end{cases}$$

$$(4)$$

Remark 2.1: If we assume that the covariance matrix for the units from level 1 in model (3) is a positive definite matrix Σ_0 instead of the identity matrix, when Σ_0 is known, the two assumptions are actually equivalent by taking reparameterization that $y_{i_0}^* =$ $\Sigma_0^{-\frac{1}{2}} y_{i_0}$ and $Z_{i_0}^* = \Sigma_0^{-\frac{1}{2}} Z_{i_0}$. Furthermore, for a technical reason, Σ_0 must be assumed as known throughout this paper, and this reason will be explained in Section 5.

2.2. Connection with the linear mixed-effect model (LMM)

The two-level hierarchical normal models are often referred to as LMMs in many places. As for the GHNL model (4), let $\Theta = \{\theta_1, \ldots, \theta_r\}$ denote the set of unobserved outcome vectors and $\mathcal{V} = \{V_1, \ldots, V_r\}$ represent the set of unknown covariance matrices. If we take θ_j 's as intermediate variables, then marginalizing out over Θ yields

$$(\boldsymbol{y}|\boldsymbol{\eta},\mathcal{V}) \sim N_{m_0k_0}(\boldsymbol{X}_r\boldsymbol{\eta},\boldsymbol{\Delta}),$$
 (5)

where

$$\Delta = I_{m_0k_0} + \sum_{t=1}^{j} X_{t-1} (I_{m_t} \otimes V_t) X_{t-1}^{\top}, \text{ and}$$
$$X_j = \prod_{s=0}^{j} Z_s, \, j = 0, 1, \dots, r.$$
(6)

\Delta is a $(m_0k_0) \times (m_0k_0)$ matrix and X_j are $(m_0k_0) \times (m_{j+1}k_{j+1})$ matrices for j = 0, 1, ..., r. Suppose that Z_j are of full column ranks. Then by Sylvester's rank inequality X_j are also of full column ranks, j = 0, 1, ..., r. In the rest of the paper, Z_j are assumed to be of full column ranks for j = 0, 1, ..., r.

If we consider a particular LMM as

$$y = X_r \eta + X_0 \theta_1^* + \dots + X_{r-1} \theta_r^* + \epsilon, \qquad (7)$$

where η is the fixed effect, θ_j 's are random effects and independently distributed as $N_{m_jk_j}(\mathbf{0}, \mathbf{I}_{m_j} \otimes \mathbf{V}_j)$ for $j \in [r]$, and $\boldsymbol{\epsilon}$ denotes the vector of random errors and is distributed as $N_{m_0k_0}(\mathbf{0}, \mathbf{I}_{m_0k_0})$. By integrating out the random effects, the marginal distribution of \boldsymbol{y} conditioning on (η, \mathcal{V}) is identical to the distribution (5). In one word, the GHNL is equivalent to an LMM in the sense of the marginal distribution of observations after integrating out intermediate outcome vectors or random effects. The equivalence between the GHNL models and the LMMs can be illustrated by an example of mixed-effect ANOVA model as well.

Example 2.1 (Mixed-effect ANOVA model): Suppose we can observe the scores of *p* courses for student (*ijk*) as y_{ijk} for $i = 1, ..., s_1, j = 1, ..., s_2$ and $k = 1, ..., s_3$. The observed data are within a hierarchy of three levels: student (*ijk*) is nested within class (*ij*), and class (*ij*) is nested within school *i*. Thus, we have total s_1 schools, each school has s_2 classes and each class has s_3 students. Consider a mixed-effect ANOVA model as

$$\boldsymbol{y}_{ijk} = \boldsymbol{\eta} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_{ij} + \boldsymbol{\epsilon}_{ijk}, \qquad (8)$$

where \mathbf{y}_{ijk} , $\boldsymbol{\eta}$, α_i , $\boldsymbol{\beta}_{ij}$ and $\boldsymbol{\epsilon}_{ijk}$ are all $p \times 1$ vectors for $i = 1, \ldots, s_1, j = 1, \ldots, s_2$ and $k = 1, \ldots, s_3, \boldsymbol{\eta}$ denotes the overall mean and is fixed effect, $\boldsymbol{\alpha}_i \sim N_p(\mathbf{0}, V_{\alpha})$ is the effect of school, $\boldsymbol{\beta}_{ij}$ is distributed as $N_p(\mathbf{0}, V_{\beta})$

and represents the effect of class, the student-level independent random error is denoted by ϵ_{ijk} and has distribution $N_p(\mathbf{0}, \Sigma_0)$, and Σ_0 is a known matrix. At the same time, $\boldsymbol{\alpha}_i$, $\boldsymbol{\beta}_{ij}$ and ϵ_{ijk} are independently distributed. Consequently, V_{α} , V_{β} and Σ_0 are the variance components describing the school-level, class-level and student-level variations, respectively. Due to the hierarchical structure of the observations, we can naturally build a hierarchical model as

$$\mathbf{y}_{ijk} \sim N_p(\boldsymbol{\beta}_{ij}^*, \boldsymbol{\Sigma}_0), \quad \boldsymbol{\beta}_{ij}^* \sim N_p(\boldsymbol{\alpha}_i^*, \boldsymbol{V}_\beta) \quad \text{and}$$

 $\boldsymbol{\alpha}_i^* \sim N_p(\boldsymbol{\eta}, \boldsymbol{V}_{\boldsymbol{\alpha}})$ (9)

independently, for $i = 1, ..., s_1$, $j = 1, ..., s_2$ and $k = 1, ..., s_3$. Denote that

$$Y_{i} = \begin{pmatrix} y_{i11} & \cdots & y_{is_{2}1} \\ \vdots & \ddots & \vdots \\ y_{i1s_{3}} & \cdots & y_{is_{2}s_{3}} \end{pmatrix},$$
$$E_{i} = \begin{pmatrix} \epsilon_{i11} & \cdots & \epsilon_{is_{2}1} \\ \vdots & \ddots & \vdots \\ \epsilon_{i1s_{3}} & \cdots & \epsilon_{is_{2}s_{3}} \end{pmatrix},$$
$$\beta_{i} = \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{is_{2}} \end{pmatrix} \text{ and }$$
$$\beta_{i}^{*} = \begin{pmatrix} \beta_{i1} \\ \vdots \\ \beta_{is_{2}} \end{pmatrix},$$

where Y_i and E_i are both $(s_3p) \times s_2$ matrices, and β_i and β_i^* are both $(s_2p) \times 1$ vectors. Let $y_i = \text{vec}(Y_i)$ and $\epsilon_i = \text{vec}(E_i)$, where vec(A) denotes the column vector obtained by stacking the columns of the matrix A on top of one another. Define that

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_{s_1} \end{pmatrix}, \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_{s_1} \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_{s_1} \end{pmatrix},$$
$$\beta^* = \begin{pmatrix} \beta_1^* \\ \vdots \\ \beta_{s_1}^* \end{pmatrix}, \quad \alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_{s_1} \end{pmatrix}, \quad \alpha^* = \begin{pmatrix} \alpha_1^* \\ \vdots \\ \alpha_{s_1}^* \end{pmatrix}.$$

Thus, \boldsymbol{y} and $\boldsymbol{\epsilon}$ are $(m_0 p) \times 1$ vectors, and $\boldsymbol{\beta}$ and $\boldsymbol{\beta}^*$ are $(m_1 p) \times 1$ vectors, and $\boldsymbol{\alpha}$ and $\boldsymbol{\alpha}^*$ are $(m_2 p) \times 1$ vectors, where $m_0 = s_1 s_2 s_3$, $m_1 = s_1 s_2$ and $m_2 = s_1$. Then the hierarchical normal model (9) can be expressed as a GHNL model with the form

$$\begin{cases} \text{Level 1:} \quad (\boldsymbol{y}|\boldsymbol{\beta}^{*}, \boldsymbol{\Sigma}_{0}) & \sim & N_{m_{0}p}(\boldsymbol{Z}_{0}\boldsymbol{\beta}^{*}, \boldsymbol{I}_{m_{0}} \otimes \boldsymbol{\Sigma}_{0}); \\ \text{Level 2:} \quad (\boldsymbol{\beta}^{*}|\boldsymbol{\alpha}^{*}, \boldsymbol{V}_{\beta}) & \sim & N_{m_{1}p}(\boldsymbol{Z}_{1}\boldsymbol{\alpha}^{*}, \boldsymbol{I}_{m_{1}} \otimes \boldsymbol{V}_{\beta}); \\ \text{Level 3:} \quad (\boldsymbol{\alpha}^{*}|\boldsymbol{\eta}, \boldsymbol{V}_{\alpha}) & \sim & N_{m_{2}p}(\boldsymbol{Z}_{2}\boldsymbol{\eta}, \boldsymbol{I}_{m_{2}} \otimes \boldsymbol{V}_{\alpha}), \end{cases}$$
(10)

where

$$Z_0 = \operatorname{diag}\{\underbrace{\mathbf{1}_{s_3} \otimes \mathbf{I}_p, \dots, \mathbf{1}_{s_3} \otimes \mathbf{I}_p}_{s_1 s_2}\},$$
$$Z_1 = \operatorname{diag}\{\underbrace{\mathbf{1}_{s_2} \otimes \mathbf{I}_p, \dots, \mathbf{1}_{s_2} \otimes \mathbf{I}_p}_{s_1}\}, Z_2 = \mathbf{1}_{s_1} \otimes \mathbf{I}_p,$$

where $\mathbf{1}_q$ denotes the $q \times 1$ vector with all elements being one, and \mathbf{Z}_0 , \mathbf{Z}_1 , \mathbf{Z}_2 are $(m_0p) \times (m_1p)$, $(m_1p) \times (m_2p)$, $(m_2p) \times p$ matrices, respectively. Denote that

$$X_0 \stackrel{\Delta}{=} Z_0,$$

$$X_1 \stackrel{\Delta}{=} \operatorname{diag}\{\underbrace{\mathbf{1}_{s_2 s_3} \otimes I_p, \dots, \mathbf{1}_{s_2 s_3} \otimes I_p}_{s_1}\} = Z_0 Z_1,$$

$$X_2 \stackrel{\Delta}{=} \mathbf{1}_{s_1 s_2 s_3} \otimes I_p = Z_0 Z_1 Z_2,$$

and X_0 , X_1 , X_2 are $(m_0p) \times (m_1p)$, $(m_0p) \times (m_2p)$, $(m_0p) \times p$ matrices, respectively. Thus, model (8) can be summarized as

$$y = X_2 \eta + X_1 \alpha + X_0 \beta + \epsilon, \qquad (11)$$

where $\boldsymbol{\alpha} \sim N_{m_2p}(\mathbf{0}, \boldsymbol{I}_{m_2} \otimes \boldsymbol{V}_{\alpha}), \quad \boldsymbol{\beta} \sim N_{m_1p}(\mathbf{0}, \boldsymbol{I}_{m_1} \otimes \boldsymbol{V}_{\beta})$ and $\boldsymbol{\epsilon} \sim N_{m_0p}(\mathbf{0}, \boldsymbol{I}_{m_0} \otimes \boldsymbol{\Sigma}_0)$, independently. By integrating out $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ and $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$, the marginal distributions of \boldsymbol{y} for model (11) and model (10) are identical and of the form:

$$(\mathbf{y}|\boldsymbol{\eta}, \mathbf{V}_{\alpha}, \mathbf{V}_{\beta}, \boldsymbol{\Sigma}_{0}) \sim N_{m_{0}p} (\mathbf{X}_{2}\boldsymbol{\eta}, \boldsymbol{\Omega}),$$

 $\boldsymbol{\Omega} = \mathbf{I}_{m_{0}} \otimes \boldsymbol{\Sigma}_{0} + \mathbf{X}_{0} (\mathbf{I}_{m_{1}} \otimes \mathbf{V}_{\beta}) \mathbf{X}_{0}^{\top} + \mathbf{X}_{1} (\mathbf{I}_{m_{2}} \otimes \mathbf{V}_{\alpha}) \mathbf{X}_{1}^{\top}.$

Example 2.1 provides a simple example to illustrate how the hierarchical model and the mixed-effect model can be constructed based on the nested data, and the equivalence between two models is also presented. In Appendix 1, we define a special LMM which is a special case of model (7) with $m_j = 1$ for all $j \in [r]$, and the theoretical investigation of this special LMM is distinct from that of the GHNLM. This special LMM could be common in application, and we just want to provide some theoretical results for those who have interests to refer to. We will still focus on the GHNL model in the following sections.

2.3. Priors on the hyperparameters

In order to implement fully Bayesian analysis, we should specify hyperpriors on the parameters (η , \mathcal{V}). It follows the recommendation from Berger et al. (2020b) that we can assume priors on (η , \mathcal{V}) as:

$$\pi(\eta) \propto \frac{1}{(1+\|\eta\|^2)^{(d-1)/2}}, \quad \eta \in \mathbb{R}^d,$$
(12)
$$\pi(V_j) \propto \frac{1}{|V_j|^{1-1/(2k_j)} \prod_{s < t} (\omega_{js} - \omega_{jt})},$$

$$V_j > 0, \ j \in [r],$$
 (13)

where $\omega_{j1} > \omega_{j2} > \cdots > \omega_{jk_j} > 0$ are the decreasingly ordered eigenvalues of V_j , $j \in [r]$. Apart from prior (12), common choices of prior on η include the constant prior and conjugate prior. None of the three priors will result in improper priors or difficulties in computation. However, among the three priors, prior (12) is the most perfect for all k from the perspective of admissibility. Besides, it refers to Berger et al. (2005) that the prior (12) is a mixture-of-normal prior of the hierarchical structure as

$$(\boldsymbol{\eta}|\boldsymbol{\lambda}) \sim N_d(\mathbf{0}, \boldsymbol{\lambda} \boldsymbol{I}_d) \text{ and } [\boldsymbol{\lambda}] \propto \boldsymbol{\lambda}^{-1/2} \exp\left(-\frac{1}{2\boldsymbol{\lambda}}\right),$$

(14)

and those mixture-of-normal priors have shown great success in shrinkage estimation particularly (cf. Fourdrinier et al., 1998) and robust Bayesian estimation generally (cf. Berger, 1980). Therefore, prior (12) was actually recommended by Berger et al. (2005) for default use.

As for prior (13) on the unknown covariance matrices $V_j, j \in [r]$, consider the transformation from V_j to $\Omega_j = \text{diag}(\omega_{j1}, \ldots, \omega_{jk_j})$ and the orthogonal matrix Γ_j of corresponding eigenvectors, the Jacobian is

$$\left|\frac{\partial \mathbf{V}_j}{\partial \left(\mathbf{\Omega}_j, \mathbf{\Gamma}_j\right)}\right| = \prod_{s < t} \left(\omega_{js} - \omega_{jt}\right). \tag{15}$$

Consequently, the prior (13) on V_j becomes the prior density of (Ω_j, Γ_j) as

$$\pi(\mathbf{\Omega}_j, \mathbf{\Gamma}_j) \propto \frac{1}{|\mathbf{\Omega}_j|^{1-1/(2k_j)}}$$
(16)

with respect to Lebesgue measure on $(\omega_{j1}, \ldots, \omega_{jk_j})$ and the invariant Haar measure over the space $\{\boldsymbol{\Gamma}: \boldsymbol{\Gamma}\boldsymbol{\Gamma}^{\top} = \boldsymbol{I}_{k_j}\}$. Note that the prior on $\boldsymbol{\Omega}_j$ is improper and, independently, the prior on Γ_i is constant. Use of a uniform prior for Γ_i ranging over a compact space is natural and non-controversial and has no influence on the eigenvalues. The term $\prod_{s < t} (\omega_{js} - \omega_{jt})$ is eliminated after changing variables for prior (13). In contrast, the commonly used priors on the covariance matrix, such as inverse Wishart, Jeffreys-rule and constant priors, contain the term $\prod_{s < t} (\omega_{js} - \omega_{jt})$ in the transformed space. This special term gives low mass to close eigenvalues and hence effectively force the eigenvalues apart. It is contrary to the common intuition which would suggest that one should choose a prior that pushes the eigenvalues closer together. As a result, prior (13) is essentially neutral as to expansion or shrinkage of the eigenvalues.

In the context of 2-level hierarchical normal model, Theorem 1 from Berger et al. (2020b) has demonstrated that the combination of priors (12) and (13) on (η , V) is on the boundary of admissibility, being as diffuse as possible without yielding inadmissible estimators. Furthermore, it is shown that the generalization that allows covariates at all levels of the hierarchical model will not affect the result of admissibility (cf. Berger et al., 2020b). Nonetheless, the admissibility of the recommended prior for the (r + 1)-level hierarchical model with $r \ge 2$ is not clear. Generally speaking, this is a very difficult question to answer, and we mainly justify the recommendation of hyperpriors from the angles of posterior propriety and computation afterwards in the framework of the GHNL model.

3. Posterior propriety

Berger et al. (2020b) has shown that the resulting posterior of the recommended prior is proper for the 3level hierarchical model (2), but under a narrow set of assumptions. They also conjectured the posterior propriety for a hierarchical model with any number of levels, a rigorous one of which was not given yet. In this section, we will comprehensively investigate the conditions for the posterior propriety of the GHNL model (4) using the recommended prior in more general situations. The dimension of η affects the investigation of posterior propriety considerably, and two cases $d \ge 2$ or d = 1 will be discussed separately.

Based on (5) and (14), by integrating out η , we can obtain the marginal distribution of *y* conditioning on (\mathcal{V}, λ) as

$$(\boldsymbol{y}|\mathcal{V},\lambda) \sim N_{m_0k_0}\left(\boldsymbol{0},\boldsymbol{\Delta}+\lambda\boldsymbol{X}_r\boldsymbol{X}_r^{\top}\right).$$
 (17)

The posterior propriety of the GHNL model (4) employing priors (13) and (14) is defined as

$$m(\mathbf{y}) = \int f(\mathbf{y}|\mathcal{V}, \lambda) \pi(\lambda) \prod_{s=1}^{r} \pi(\mathbf{V}_{s}) \,\mathrm{d}\lambda \prod_{t=1}^{r} \,\mathrm{d}\mathbf{V}_{t} < \infty,$$
(18)

where m(y) denotes the marginal density of the observation vector. Next, we display some definitions and additional notations which are frequently used in this section.

More Notations Let card(A) denote the cardinality of the set A. For $0 < I_1, I_2 \le \infty$, denote that $I_1 \simeq I_2$ if there exist constants $0 < C_1 \le C_2$ such that $C_1I_2 \le I_1 \le C_2I_2$; For a symmetric matrix $A \in \mathbb{R}^{n \times n}$, $\lambda_i(A)$ represents the *i*-th largest eigenvalue of A, namely, $\lambda_1(A) \ge \cdots \ge \lambda_n(A)$; Let $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ denote the maximum and minimum eigenvalues of an arbitrary symmetric matrix A.

Definition: For convenience, let $\omega_{r+1,1} \triangleq \lambda$, $m_{r+1} \triangleq d$ and $k_{r+1} \triangleq 1$. Define a function of an arbitrary nonempty set *E* as $\mathcal{F}(E) \triangleq \{D \mid D \subset E, D \neq \emptyset\}$, such that $\mathcal{F}(E)$ denotes the set of all non-empty subsets of *E*. For any $R \in \mathcal{F}([r+1])$, let $H(R) \triangleq \{(j,l) | j \in R, l \in [k_j]\}$. Further define the composition of \mathcal{F} and H as

$$\mathcal{S}(R) \triangleq (\mathcal{F} \circ H)(R) = \{D \mid D \subset H(R), \ D \neq \emptyset\}$$

for any $R \in \mathcal{F}([r+1])$. Define that

$$c_{jl,s} = 1_{\{m_j(l-1) < s \le m_j l\}},\tag{19}$$

for $j \in [r + 1]$, $l \in [k_j]$ and $s \in [m_0k_0]$.

3.1. Two key lemmas

Before we formally investigate the posterior propriety, we first introduce two important lemmas which dominate in the process of proving the main theorems in this paper.

Lemma 3.1: Assume that A_j are $p_j \times p_j$ positive definite matrices, $j \in [r]$. Let X_j be $n \times p_j$ matrices of full column ranks, $j \in [r]$. Define that

$$H = \sum_{j=1}^{r} X_j A_j X_j^{\top}.$$
 (20)

Then

(a) $\lambda_{\max}(H) \leq C_1 \sum_{j=1}^r \lambda_{\max}(A_j)$, where $C_1 = \max_{j \in [r]} \lambda_{\max}(X_j^\top X_j)$.

(b) Also,

$$|\boldsymbol{H}| \geq \left(\frac{C_2}{r}\right)^n \left|\sum_{j=1}^r \boldsymbol{D}_j\right|, \qquad (21)$$

where $C_2 = \min_{j \in [r]} \lambda_{\min}(\mathbf{X}_j^{\top} \mathbf{X}_j) > 0$. For any $j \in [r]$, $\mathbf{D}_j = \operatorname{diag}(a_{j1}, \ldots, a_{jn})$, where $a_{jk} = \lambda_k (\mathbf{A}_j)$ for $k \in [p_j]$ and $a_{jk} = 0$ for $p_j < k \leq n$.

Lemma 3.1 mainly demonstrates two inequalities with respect to the summation of a series of quadratic forms which have matrical inputs. Since X_i 's have full column rank and A_j 's are positive definite, then $n \ge p_j$ and rank $\left(X_j A_j X_j^{\top}\right) = p_j$. It is worthwhile to note that the non-zero diagonal elements of D_j are the decreasingly ordered eigenvalues of A_j in the lower bound of |H|, and this relation will deeply influence the last result when we derive the sufficient condition for posterior propriety afterwards. Besides, we can never find a constant $C^* > 0$ such that $|H|_+ \le C^* \left|\sum_{j=1}^r D_j\right|_+$, where $|M|_+$ denotes the product of all non-zero eigenvalues of M. The proof of Lemma 3.1 can be found in Appendix A.2.

Lemma 3.2: Assume a positive integer k. Suppose \mathcal{M} is a subset of $\mathcal{F}([k])$ with cardinality n and let C_1, \ldots, C_n

denote the sequence of all the elements within \mathcal{M} . Define an integral

$$I = \int_{\Omega} \frac{1}{\left[\prod_{j=1}^{k} \lambda_j^{a_j}\right] \left[\prod_{i=1}^{n} (1 + \sum_{r \in C_i} \lambda_r)^{b_i}\right]} d\lambda_i$$

where $a_j \ j = 1, 2, ..., k$, and b_i , i = 1, 2, ..., n, are real constants, and $\lambda = (\lambda_1, ..., \lambda_k)^\top \in \Omega \equiv [0, \infty]^k$. Then the integral I is finite if and only if (iff) the following two conditions are both satisfied.

(a) *a_j* < 1, *j* ∈ [*k*];
(b) *inequalities*

$$\sum_{j\in D} a_j + \sum_{i\in \mathcal{G}_D} b_i > \operatorname{card}(D)$$
(22)

for all $D \in \mathcal{F}([k])$ hold, where $\mathcal{G}_D = \{i \mid D \bigcap C_i \neq \emptyset, i \in [n]\}.$

Here, Lemma 3.2 (b) may not be straightforward enough, so we will use the following example to elaborate on how Lemma 3.2 can be employed.

Example 3.1: Consider integral

$$I_0 = \int_{\Omega} \frac{1}{\left[\prod_{j=1}^3 \lambda_j^{a_j}\right] (1+\lambda_1)^{b_1} (1+\lambda_1+\lambda_2)^{b_2} (1+\lambda_2+\lambda_3)^{b_3}} \, d\lambda_1 \, d\lambda_2 \, d\lambda_3,$$

where a_1 , a_2 , a_3 , b_1 , b_2 , b_3 are all real constants. Similar to Lemma 3.2, we can define $\mathcal{M} = \{\{1\}, \{1, 2\}, \{2, 3\}\}$. Then $I_0 < \infty$ iff all the following inequalities hold: (a) $a_j < 1$, $j \in [3]$; (b)

for
$$D = \{1\}$$
, $a_1 + b_1 + b_2 > 1$;
for $D = \{2\}$, $a_2 + b_2 + b_3 > 1$;
for $D = \{3\}$, $a_3 + b_3 > 1$;
for $D = \{1, 2\}$, $a_1 + a_2 + b_1 + b_2 + b_3 > 2$;
for $D = \{1, 3\}$, $a_1 + a_3 + b_1 + b_2 + b_3 > 2$;
for $D = \{2, 3\}$, $a_2 + a_3 + b_2 + b_3 > 2$;
for $D = \{1, 2, 3\}$, $a_1 + a_2 + a_3 + b_1 + b_2 + b_3 > 3$.

Note that, no matter how \mathcal{M} is defined, we always need to check all the inequalities corresponding to all $D \in \mathcal{F}([k])$. Even though some inequalities could be trivial after being written down, for the sake of assurance, we would better take all non-empty subsets of [k] into account in the early stage. Lemma 3.2 plays a crucial role in obtaining the follow-up theorems, and detailed proof of Lemma 3.2 sees Appendix A.2.

3.2. Conditions for the posterior to be proper when $d \ge 2$

In this subsection, the case $d \ge 2$ is mainly considered.

Theorem 3.1: Consider the GHNL model (4) with priors (12) and (13) on η and $V_i \in \mathcal{V}$, respectively. When

 $d \ge 2$, a sufficient condition for the posterior propriety is given by

$$\sum_{s=1}^{m_0k_0} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > \sum_{(j,l)\in D} \frac{1}{k_j},$$
 (23)

for any $D \in \mathcal{S}([r+1])$.

Proof: It follows (17) that the integrated likelihood of (\mathcal{V}, λ) after marginalizing out $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r, \eta)$ is given by

$$L(\mathcal{V}, \lambda; \mathbf{y}) \propto \frac{1}{|\mathbf{\Delta} + \lambda \mathbf{X}_r \mathbf{X}_r^{\top}|^{1/2}} \times \exp\left\{-\frac{1}{2}\mathbf{y}^{\top} \left(\mathbf{\Delta} + \lambda \mathbf{X}_r \mathbf{X}_r^{\top}\right)^{-1} \mathbf{y}\right\}.$$
(24)

By dropping the exponent term involving y (since it is less than one), we have

$$L(\mathcal{V}, \lambda; \mathbf{y}) < \frac{1}{|\mathbf{\Delta} + \lambda \mathbf{X}_r \mathbf{X}_r^{\top}|^{1/2}}$$

By applying Lemma 3.1 (b), we can further bound the integrated likelihood as

$$L(\mathcal{V},\lambda;\mathbf{y}) \le \frac{C_1}{|\mathbf{M}_1|^{1/2}} \tag{25}$$

where C_1 is a positive constant that only depends on X_j 's and

$$M_{1} = I_{m_{0}k_{0}} + \sum_{j=1}^{r+1} D_{j} \text{ and}$$
$$D_{j} = \begin{pmatrix} \mathbf{\Omega}_{j} \otimes I_{m_{j}} \\ & \mathbf{O}_{q_{j}} \end{pmatrix}_{(m_{0}k_{0}) \times (m_{0}k_{0})}, \quad (26)$$

where Ω_j are the diagonal matrices of the decreasingly ordered eigenvalues of V_j for $j \in [r]$, $\Omega_{r+1} = (\omega_{r+1,1})$, O_{q_j} are $q_j \times q_j$ zero matrices and $q_j = m_0 k_0 - m_j k_j$ for $j \in [r+1]$. Since $\Omega_{r+1} = (\omega_{r+1,1})$, then Ω_{r+1} is a degenerate matrix as scalar λ and the prior on λ becomes

$$\pi(\mathbf{\Omega}_{r+1}) \propto \frac{1}{|\mathbf{\Omega}_{r+1}|^{1-1/(2k_{r+1})}} \exp\left\{-\frac{1}{2}\operatorname{tr}\left(\mathbf{\Omega}_{r+1}^{-1}\right)\right\}.$$
(27)

Combining (16), (18), (27) and (25), we have

$$\begin{split} m(\mathbf{y}) &\leq \int \frac{C_1 \exp\left\{-\operatorname{tr}\left(\mathbf{\Omega}_{r+1}^{-1}\right)/2\right\} \prod_{j=1}^{r+1} \mathbf{1}_{\left\{\omega_{j1} > \cdots > \omega_{jk_j} > 0\right\}}}{|M_1|^{1/2} \prod_{j=1}^{r+1} |\mathbf{\Omega}_j|^{1-\frac{1}{2k_j}}} \\ &\times \left[\prod_{j=1}^r \,\mathrm{d}\mathbf{\Omega}_j \,\mathrm{d}\mathbf{\Gamma}_j\right] \,\mathrm{d}\mathbf{\Omega}_{r+1} \\ &< \int \frac{C_1}{|M_1|^{1/2} \prod_{j=1}^{r+1} |\mathbf{\Omega}_j|^{1-\frac{1}{2k_j}}} \prod_{j=1}^{r+1} \,\mathrm{d}\mathbf{\Omega}_j \triangleq I_0. \end{split}$$

The definition of $c_{jl,s}$ in (19) yields

$$|M_1| = \prod_{s=1}^{m_0 k_0} \left(1 + \sum_{j=1}^{r+1} \sum_{l=1}^{k_j} c_{jl,s} \omega_{jl} \right).$$

Therefore,

$$I_0 \simeq \int \frac{\prod_{j=1}^{r+1} \, \mathrm{d}\mathbf{\Omega}_j}{\prod_{j=1}^{r+1} \left|\mathbf{\Omega}_j\right|^{1-\frac{1}{2k_j}} \prod_{s=1}^{m_0 k_0} \left(1 + \sum_{j=1}^{r+1} \sum_{l=1}^{k_j} c_{jl,s} \omega_{jl}\right)^{\frac{1}{2}}},$$

which is finite iff

$$\sum_{(j,l)\in D} \left(1 - \frac{1}{2k_j}\right) + \frac{1}{2} \sum_{s=1}^{m_0 k_0} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > \operatorname{card}(D)$$
(28)

for any $D \in \mathcal{S}([r+1])$ by employing Lemma 3.2. It's obvious that the inequality (28) is equivalent to that in (23).

Applying Lemma 3.1 yields the upper bound on the integrated likelihood function of hyperparameters as $C_0|M_1|^{-\frac{1}{2}}$, where M_1 is an $(m_0k_0) \times (m_0k_0)$ matrix and defined in (26). Therefore, the special notation $c_{jl,s}$ can be understood as the indicator of whether eigenvalue ω_{jl} appears in the *s*-th diagonal element of M_1 for $j \in [r+1]$, $l \in [k_j]$ and $s \in [m_0k_0]$. At the same time, the left-hand side of inequality (23) actually stands for the cardinality of set $\{s \mid \exists (j, l) \in D, \text{ such that } c_{jl,s} > 0\}$ for any $D \in \mathcal{S}([r+1])$.

The cardinality of S([r+1]) is $2^{\sum_{j \in [r+1]} k_j} - 1$, which means that the total number of the inequalities to be checked in (23) is exponential with r and the dimensions k_j . It has to be admitted that this will impose considerably heavy computational burden in common practice by applying Theorem 3.1 directly. Nevertheless, the researchers have no need to be anxious about the heavy computational burden, because most of the inequalities in Theorem 3.1 are trivial. To conclude this point, we have the following corollary.

Corollary 3.1: Recursively define that

$$R_{1} = \{j \mid m_{j} \leq r+1, \quad j \in [r+1]\};$$

$$R_{2} = \{j \mid m_{j} \leq \operatorname{card}(R_{1}), \quad j \in R_{1}\};$$

$$\vdots$$

$$R_{p} = \{j \mid m_{j} \leq \operatorname{card}(R_{p-1}), \quad j \in R_{p-1}\},$$

where *p* is the smallest positive integer *i* such that $\{j | m_j > card(R_i), j \in R_i\} = \emptyset$. We call the levels within R_p as kernel levels and denote R_p by R_{ker} . The inequality (23) holds for any $D \in S([r+1])$ iff the inequality (23) holds for any $D \in S(R_{ker})$. Consequently, if $R_{ker} = \emptyset$, then the posterior is always proper.

Proof: Let $R_1^c = [r+1] - R_1$, thus, $m_j > r+1$ for $j \in R_1^c$ by the definition of R_1 . For any $D \in S([r+1])$, if there exists $j^* \in R_1^c$ such that $(j^*, l) \in D$ for any $l \in [k_j^*]$, then

$$\sum_{s=1}^{m_0 k_0} \mathbb{1}_{\left\{\sum_{(j,l) \in D} c_{jl,s} > 0\right\}} \ge \max_{(j,l) \in D} m_j > r+1$$
$$= \max_{D \in \mathcal{S}([r+1])} \sum_{(j,l) \in D} \frac{1}{k_j}$$

which is a trivial one. As a result, inequality (23) holds for any $D \in S([r + 1])$ iff inequality (23) holds for any $D \in S(R_1)$. Since

$$\max_{D\in\mathcal{S}(R_i)}\sum_{(j,l)\in D}\frac{1}{k_j}=\operatorname{card}(R_i),$$

it can be recursively shown that inequality (23) holds for any $D \in S([r + 1])$ iff inequality (23) holds for any $D \in S(R_i)$ and $i \in [p]$, where *p* is the smallest positive integer *i* such that $R_i - R_{i+1} = \emptyset$.

By using the technique of extracting kernel levels, we dramatically narrow down the checking region for posterior propriety. Since we should only check the inequalities for the levels within R_{ker} , substantially reducing the number of inequalities to be checked from $2^{\sum_{j \in [r+1]} k_j} - 1$ to $2^{\sum_{j \in R_{\text{ker}} k_j}} - 1$. Moreover, Corollary 3.1 also indicates two interesting conclusions depicted as follows.

(a) First, it reveals the mechanism how three roles, number of levels, numbers of units in levels and dimensions of levels, affect the posterior propriety simultaneously. Roughly speaking, in the context of GHNL, the more levels with fewer units having lower dimensions, the less likely the posterior is to be proper. For example, if m_{r-2} = m_{r-1} = m_r = 2 and k_{r-2} = k_{r-1} = k_r = 1, for D = {(j, l) | j = r - 2, r - 1, r, l = 1}, we have

$$2 = \sum_{s=1}^{m_0 k_0} \mathbb{1}_{\left\{\sum_{(j,l) \in D} c_{jl,s} > 0\right\}} < \sum_{(j,l) \in D} \frac{1}{k_j} = 3.$$
(29)

Therefore, the posterior propriety can hardly be guaranteed by Theorem 3.1. Conversely, if the units in each level are adequate enough such that the set of kernel levels is empty, i.e., $R_{ker} = \emptyset$, then the posterior is always proper. As a consequence, more attention should be focused on the levels with small number of units, namely, the kernel levels.

(b) Second, the recommended prior for use at any level in hierarchical modelling is further justified from the aspect of posterior propriety and ease of implementation. For instance, if we switch the prior on $V_j, j \in [r]$ from (13) to

$$\pi(\mathbf{V}_j) \propto \frac{1}{|\mathbf{V}_j|^{1-a} \prod_{1 \le s < t \le k_j} (\omega_{js} - \omega_{jt})},$$
$$\mathbf{V}_j > 0, \ j \in [r], \tag{30}$$

where $0 < a \le 1$ (*a* has to be larger than zero by Lemma (3.2)). Then the condition in Theorem 3.1 becomes

$$\sum_{s=1}^{m_0k_0} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > 2a \times \text{card}(D),$$
$$D \in \mathcal{S}([r+1]). \tag{31}$$

On one hand, when *a* is greater than but close to zero (denoted by $a \gtrsim 0$), all the inequalities in (31) always hold. Thus, the posterior will be always proper. However, it is impractical to decide how small is for *a* and find such a fixed value that fits all levels. On the other hand, when $a \lesssim$ 1, which denotes that *a* is less than and close to 1, then $2a * \operatorname{card}(D) > \sum_{(j,l)\in D} \frac{1}{k_j}$, concluding that inequality (31) is harder to get reached than inequality (23), especially for large dimensions k_j 's. Therefore, the posterior using prior (30) is rather unlikely to be proper than using prior (13). Similar to Corollary 3.1, we can recursively define that

$$R_{1}^{*} = \left\{ j \mid m_{j} \leq 2a \times \text{card}(E_{0}), \quad j \in [r+1] \right\};$$

$$R_{2}^{*} = \left\{ j \mid m_{j} \leq 2a \times \text{card}(E_{1}^{*}), \quad j \in R_{1}^{*} \right\};$$

$$\vdots$$

$$R_{p^{*}}^{*} = \left\{ j \mid m_{j} \leq 2a \times \text{card}(E_{p^{*}-1}^{*}), \quad j \in R_{p^{*}-1}^{*} \right\};$$

where $E_0 = H([r+1])$, $E_i^* = H(R_i^*)$ for $i \in [p^*]$ and p^* is the smallest positive integer l such that $\{j \mid m_j > 2a * \operatorname{card}(E_l^*), j \in R_l^*\} = \emptyset$. Then the posterior using prior (30) instead is proper if (31) holds for any $D \in S(R_{p^*}^*)$. When $a \leq$ 1, $\operatorname{card}(R_{p^*}^*)$ will be remarkably larger than $\operatorname{card}(R_{\operatorname{ker}})$ for large values of k_j 's, imposing a dramatically heavier burden of checking inequalities than that for prior (13). Above all, one sensible choice for a is that let a be inversely proportional to k_j for level j, which takes both practical and theoretical considerations into account.

The upper bound on $\sum_{(j,l)\in D} \frac{1}{k_j}$ as card (R_i) for $D \in S(R_i)$, leads to an effective way to extract kernel levels as presented in Corollary 3.1, but this upper bound is still too rough. Next, an elaborate upper bound on $\sum_{(j,l)\in D} \frac{1}{k_j}$ is demonstrated and a sufficient condition of clean form for posterior propriety is derived then.

Theorem 3.2: Consider the GHNL model (4) with priors (12) and (13) on η and $V_i \in \mathcal{V}$, respectively. Denote that $m^* = \min_{j \in [r+1]} m_j = \min_{j \in R_{ker}} m_j$. When $d \ge 2$, the posterior is always proper if

$$\sum_{j \in R_{\rm ker}} \frac{1}{k_j} < m^*. \tag{32}$$

Proof: For $D \in \mathcal{S}([r+1])$, define that

$$L(D) = \sum_{s=1}^{m_0 k_0} 1_{\left\{\sum_{(j,l) \in D} c_{jl,s} > 0\right\}}, \text{ and}$$
$$R(D) = \left\{ j \mid (j,l) \in D \right\}.$$
(33)

It follows from Corollary 3.1 that we only need to prove

$$\sum_{j,l)\in D} \frac{1}{k_j} < L(D), \quad \forall D \in \mathcal{S}(R_{\text{ker}})$$

Also, for any *D* belonging to $S(R_{ker})$, we have

$$\sum_{(j,l)\in D} \frac{1}{k_j} = \sum_{j\in R(D)} \frac{1}{k_j} \sum_{l=1}^{k_j} \mathbb{1}_{\{(j,l)\in D\}}$$
$$\leq \left(\sum_{j\in R(D)} \frac{1}{k_j}\right) \max_{j\in R(D)} \sum_{l=1}^{k_j} \mathbb{1}_{\{(j,l)\in D\}}.$$
 (34)

Distinct eigenvalues from the same level ($\leq r$ -th) never occur in the same row of matrix M_1 . Mathematically, $c_{jl_1,s}c_{jl_2,s} = 0$, for $1 \leq l_1 < l_2 \leq k_j$, $j \in [r]$, $s \in [m_0k_0]$. Thus, for any $j \in [r]$,

$$\sum_{l=1}^{k_j} \mathbb{1}_{\{(j,l)\in D\}} \le \frac{1}{m_j} L(D), \tag{35}$$

which is also true for j = r + 1 since $k_{r+1} = 1$. It is obvious that $\min_{j \in [r+1]} m_j = \min_{j \in R_{ker}} m_j$ by the definition of R_{ker} , and that is denoted by m^* . Combining (35) with (34) yields

$$\sum_{(j,l)\in D} \frac{1}{k_j} \le \left(\sum_{j\in R(D)} \frac{1}{k_j}\right) \frac{L(D)}{\min_{j\in R(D)} m_j}$$
$$\le \frac{1}{m^*} \left(\sum_{j\in R_{ker}} \frac{1}{k_j}\right) L(D) < L(D).$$
(36)

According to the proof procedure of Theorem 3.2, it can be deduced that $\sum_{j \in [r+1]} \frac{1}{k_j} < m^*$ is also sufficient for posterior propriety. Obviously, the condition in Theorem 3.2 is easier to be satisfied. Theorem 3.2 reveals that for fixed m^* , the posterior is more likely to be proper for higher dimensions of the units in the kernel levels. Theorem 3.2 also provides the researchers

with a powerful tool to check the posterior propriety quickly.

Remark 3.1: Consider the model (4) with r = 1, namely, a two-level hierarchical model. When $d \ge 2$, we have $m^* = \min \{d, m_1\} \ge 2$. Then the posterior using the recommended prior is always proper for $k_1 \ge 2$ referring to Theorem 3.2.

Example 3.2 (Continue with Example 2.1): Consider the GHNL modelling of the mixed-effect ANOVA as (10), which is a 3-level hierarchical model with r = 2, $m_0 = s_1s_2s_3$, $m_1 = s_1s_2$, $m_2 = s_1$, $m_3 = p$, $k_0 = k_1 = k_2 = p$ and $k_3 = 1$. It is natural to assume that we have at least two schools, each school has at least two classes and each class has at least two students, namely, $s_1 \ge 2$, $s_2 \ge 2$ and $s_3 \ge 2$. If (a) p > 2 and $s_1 \ge 2$ or (b) $p \ge 2$ and $s_1 > 2$ holds, it can be readily derived that the set of kernel levels is empty. Thus, the posterior is always proper according to Corollary 3.1. When $s_1 = p = 2$, the set of kernel levels is $R_2 = \{2, 3\}$, since $\frac{1}{k_2} + \frac{1}{k_3} < 2$, then the posterior is proper by applying Theorem 3.2. In conclusion, the posterior using the recommended prior is always proper when $p \ge 2$.

In Berger et al. (2020b)'s work, for a technical reason, they assumed k = sp for 3-level hierarchical normal model (2) such that the design matrices for units within level 2 and 3 are square matrices. They eventually reached a conclusion that the posterior employing the recommended prior in model (2) is always proper for k = sp and $p \ge 2$ (p is the dimension of hypermean in model (2)). However, the assumption that the design matrices appears to be unnatural and hard to be interpreted in practice. Nevertheless, we still generalize Berger et al. (2020b)'s result to the GHNL model in the following so as to draw a consistent conclusion with theirs.

Corollary 3.2: Consider the GHNL model (4) with priors (12) and (13) on η and $\mathbf{V}_i \in \mathcal{V}$, respectively. Assume that $d \ge 2$ and $k_j = m_{j+1}k_{j+1}$, $j \in [r]$. Then the posterior is always proper.

Proof: Since $m_j \ge 2$, $j \in [r+1]$, then $m^* \ge 2$. It remains to show that $\sum_{j \in R_{ker}} \frac{1}{k_j} < 2$ by Theorem 3.2. By utilizing the condition that $k_j = m_{j+1}k_{j+1}$ for $j \in [r]$ and $k_{r+1} = 1$, we have $k_j = \prod_{s=j+1}^{r+1} m_s \ge 2^{r+1-j}$ for $j \in [r]$. Thus,

$$\sum_{j \in R_{ker}} \frac{1}{k_j} \le \sum_{j \in [r+1]} \frac{1}{k_j} \le \sum_{j \in [r+1]} \frac{1}{2^{r+1-j}} = 2 - \frac{1}{2^r} < 2,$$

which completes the proof.

Above all, a general procedure for checking the posterior propriety of the GHNL models (4) employing the recommended prior for $d \ge 2$ can be summarized as follows.

Guidance for checking the posterior propriety when $d \ge 2$:

- (a) If the design matrices for each unit in each level are square matrices, then the posterior is proper, otherwise, turn to (b).
- (b) Derive the set of kernel levels, R_{ker} . If $R_{ker} = \emptyset$ or inequality (32) holds, the the posterior is proper. If neither, turn to (c).
- (c) Check inequality (23) for all *D* belonging to $S(R_{ker})$. If that always holds, the the posterior is proper. If not, the posterior propriety can hardly be guaranteed.

3.3. Conditions for the posterior to be proper when d = 1

It is quite common that the dimension of the fixed effect η is only one in practice. However, when d = 1, note that

$$1 = m_{r+1} = \sum_{s=1}^{m_0 k_0} 1_{\left\{\sum_{(j,l) \in D} c_{jl,s} > 0\right\}}$$
$$= \sum_{(j,l) \in D} \frac{1}{k_j} = \frac{1}{k_{r+1}} = 1$$

for $D = \{(r + 1, 1)\}$, resulting in the failure of the sufficient condition in Theorem 3.1. Therefore, in this subsection, we mainly reinvestigate the conditions for the the posterior to be proper for d = 1.

Theorem 3.3: Consider the GHNL model (4) with constant prior on η and prior (13) on $V_j \in \mathcal{V}$. When d = 1, the posterior is proper if

$$\sum_{s=1}^{m_0k_0} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > \mathbb{1}_{\left\{\exists j\in[r], (j,1)\in D\right\}} + \sum_{(j,l)\in D} \frac{1}{k_j}$$
(37)
holds for all $D \in \mathcal{S}([r])$.

Proof: When d = 1, vector η will degenerate into a scalar η . Hence the prior (12) on η becomes a constant prior on η . Based on (5), by integrating out over η , we can get the upper bound on the integrated likelihood of \mathcal{V} after dropping the exponential term (less than one)

$$L(\mathcal{V}) < \frac{1}{|\mathbf{\Delta}|^{\frac{1}{2}} |\mathbf{X}_r^\top \mathbf{\Delta}^{-1} \mathbf{X}_r|^{\frac{1}{2}}}.$$

$$L(\mathcal{V}) < C_0 |\mathbf{\Delta}|^{-\frac{1}{2}} \left(1 + \sum_{j=1}^r \omega_{j1} \right)^{\frac{1}{2}}$$
$$\leq C_1 \left(1 + \sum_{j=1}^r \omega_{j1} \right)^{\frac{1}{2}} |\mathbf{M}_2|^{-\frac{1}{2}}$$

where C_0 and C_1 are constants that are independent of the Ω_j for $j \in [r]$, $M_2 = I_{m_0k_0} + \sum_{j=1}^r D_j$ and D_j 's are defined in (26). Similar to the proof of Theorem 3.1, we can derive the upper bound on $m(\mathbf{y})$ as

$$m(\mathbf{y}) \leq \int \frac{C_1 \left(1 + \sum_{j=1}^r \omega_{j1}\right)^{\frac{1}{2}}}{|\mathbf{M}_2|^{1/2} \prod_{j=1}^r |\mathbf{\Omega}_j|^{1-\frac{1}{2k_j}}} \prod_{j=1}^r \mathrm{d}\mathbf{\Omega}_j \triangleq I_1.$$

It follows from the definition of $c_{jl,s}$ in (19) that

$$|M_2| = \prod_{s=1}^{m_0 k_0} \left(1 + \sum_{j=1}^r \sum_{l=1}^{k_j} c_{jl,s} \omega_{jl} \right).$$

Thus,

$$I_0 \simeq \int \frac{\left(1 + \sum_{j=1}^r \omega_{j1}\right)^{\frac{1}{2}} \prod_{j=1}^r d\mathbf{\Omega}_j}{\prod_{j=1}^r |\mathbf{\Omega}_j|^{1 - \frac{1}{2k_j}} \prod_{s=1}^{m_0 k_0} \left(1 + \sum_{j=1}^r \sum_{l=1}^{k_j} c_{jl,s} \omega_{jl}\right)^{\frac{1}{2}}},$$

which is finite iff

$$\sum_{(j,l)\in D} \left(1 - \frac{1}{2k_j}\right) + \frac{1}{2} \sum_{s=1}^{m_0 k_0} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} - \frac{1}{2} \mathbb{1}_{\left\{\exists j \in [r], (j,1)\in D\right\}} > \operatorname{card}(D)$$
(38)

for any $D \in S([r])$ by employing Lemma 3.2. It's obvious that inequality (38) is equivalent to that in (37).

Resembling the interpretation of Theorem 3.1, the left-hand side of inequality (37) actually denotes the cardinality of set $\{s \mid \exists (j, l) \in D, \text{ such that } c_{jl,s} > 0\}$ for any $D \in S([r])$. To reduce the burden of checking inequalities, we have the following corollary.

Corollary 3.3: Recursively define that

$$\tilde{R}_{1} = \{ j \mid m_{j} \leq r+1, \ j \in [r] \};$$

$$\tilde{R}_{2} = \{ j \mid m_{j} \leq \operatorname{card}(\tilde{R}_{1})+1, \ j \in \tilde{R}_{1} \};$$

$$\vdots$$

$$\tilde{R}_{q} = \{ j \mid m_{j} \leq \operatorname{card}(\tilde{R}_{p-1})+1, \ j \in \tilde{R}_{p-1} \}$$

where q is the smallest positive integer i such that $\{j \mid m_j > card(\tilde{R}_i) + 1, j \in \tilde{R}_i\} = \emptyset$. We call the levels

within \tilde{R}_q as kernel levels and denote \tilde{R}_q by \tilde{R}_{ker} . Inequality (37) holds for any $D \in S([r])$ iff inequality (37) holds for any $D \in S(\tilde{R}_{ker})$. Consequently, if $\tilde{R}_{ker} = \emptyset$, then the resulting posterior is always proper.

In the process of extracting kernel levels, the thresholds of m_j to split up levels are increased by one in Corollary 3.3, when compared with that in Corollary 3.1, and this is because the upper bound on the right-hand side of inequality (37) is increased by one than that of inequality (23). Except for this point, the proof of Corollary 3.1 is same as that of Corollary 3.3. For good measure, a simple tool to check the posterior propriety is shown as follows, which is a counterpart of Theorem 3.2 for d = 1.

Theorem 3.4: Consider the GHNL model (4) with constant prior on η and prior (13) on $V_j \in \mathcal{V}$. When d = 1, the posterior is always proper if

$$\sum_{j \in \tilde{R}_{ker}} \frac{1}{k_j} < m^* - 1, \tag{39}$$

where $m^* = \min_{j \in [r]} m_j$ and \tilde{R}_{ker} is the derived set of kernel levels.

Proof: For any $D \in S([r])$, define L(D) and R(D) in the same way as that in (33). According to Corollary (3.3), it suffices to show that

$$L(D) > \sum_{(j,l)\in D} \frac{1}{k_j} + 1$$

for any $D \in S(\tilde{R}_{ker})$. Similar to (36), we have

$$\begin{split} \sum_{(j,l)\in D} \frac{1}{k_j} &\leq \frac{1}{m^*} \left(\sum_{j\in R(D)} \frac{1}{k_j} \right) L(D) \\ &\leq \frac{1}{m^*} \left(\sum_{j\in \tilde{R}_{ker}} \frac{1}{k_j} \right) L(D) < L(D) - \frac{L(D)}{m^*}, \end{split}$$

for any $D \in S(\tilde{R}_{ker})$. Since $L(D) \ge m^*$ always holds, then the proof is completed.

Remark 3.2: In model (4), when r = 1 and d = 1, suppose $m_1 \ge 2$ and $k_1 \ge 2$. The posterior using the recommended prior is always proper by applying Theorem 3.4.

Remark 3.3: If all k_j 's are equal to one, the sufficient condition in Theorem 3.3 can be simplified as

$$\operatorname{card}(D) < \max_{j \in D} m_j - 1,$$

where $D \in \mathcal{F}([r])$. By employing the technique of extracting kernel levels, Theorem 3.4 is equivalent to Theorem 3.3, rather than a mere sufficient condition.

Example 3.3 (Continue with Example 2.1): Consider model (10) with assuming that $s_1 \ge 2$, $s_2 \ge 2$ and $s_3 \ge 2$. If p = 1, we have $k_0 = k_1 = k_2 = k_3 = m_3 = 1$. When $s_1 > 2$, we can easily derive the set of kernel levels as empty set. Thus, the posterior is always proper by Corollary 3.3. If $s_1 = 2$, inequality (37) fails, and the posterior propriety can hardly be guaranteed. Consequently, when p = 1, the posterior using the recommended prior is always proper for $s_1 > 2$.

Next, we generalize Berger et al. (2020b)'s result to the GHNL models for d = 1, assuming that the design matrices Z_{jij} 's for units are square matrices.

Corollary 3.4: When d = 1, consider the same model and prior as Theorem 3.3. Suppose $m_j \ge 3$ and $k_j = m_{j+1}k_{j+1}$, $j \in [r]$. Then the posterior is always proper.

Proof: It follows from Theorem 3.4 that we should only present

$$\sum_{j \in [r]} \frac{1}{k_j} < m^* - 1$$

According to the conditions that $k_j = m_{j+1}k_{j+1}, j \in [r]$ and $m_{r+1} = k_{r+1} = 1$, we have $k_j = \prod_{s=j+1}^r m_s, j \in [r]$. Thus,

$$\sum_{j \in [r]} \frac{1}{k_j} \le \sum_{j \in [r]} \frac{1}{3^{r-j}} = \frac{3}{2} \left(1 - \frac{1}{3^r} \right) < \frac{3}{2} < 2 \le m^* - 1.$$

Summing up the theoretical results above, a general procedure for checking the posterior propriety of the GHNL models (4) employing the recommended prior for d = 1 can be depicted as follows.

Guidance for checking the posterior propriety when d = 1:

- (a) If the design matrices for each unit in each level are square matrices and $m_j \ge 3$, $j \in [r]$, then the posterior is proper, otherwise, turn to (b).
- (b) Derive the set of kernel levels R_{ker}. If R_{ker} = Ø or inequality (39) holds, then the posterior is proper. If neither, turn to (c).
- (c) Check inequality (37) for all *D* belonging to $S(\tilde{R}_{ker})$. If that always holds, then the posterior is proper. If not, the posterior propriety can hardly be guaranteed.

4. Computation

In this section, we consider the MCMC sampling from the posterior arising from the model in Section 2. For the GHNL model (4) with prior (14) and (13) on η and \mathcal{V} , respectively, the joint posterior of $(\Theta, \mathcal{V}, \eta, \lambda)$ can be written as

$$\pi \left(\boldsymbol{\Theta}, \mathcal{V}, \boldsymbol{\eta}, \boldsymbol{\lambda} | \boldsymbol{y} \right) \propto f(\boldsymbol{y} | \boldsymbol{\theta}_1) \prod_{j=1}^{r-1} f(\boldsymbol{\theta}_j | \boldsymbol{\theta}_{j+1}, \boldsymbol{V}_j) f(\boldsymbol{\theta}_r | \boldsymbol{\eta}, \boldsymbol{V}_r)$$
$$\times \prod_{s=1}^r \pi(\boldsymbol{V}_s) \pi(\boldsymbol{\eta} | \boldsymbol{\lambda}) \pi(\boldsymbol{\lambda}). \tag{40}$$

Sampling $(\Theta, \mathcal{V}, \eta, \lambda)$ from the posterior density (40) can be handled by Gibbs sampling method. The main difficulty of the computation is to sample the covariance matrices V_i 's efficiently.

4.1. Gibbs sampling for input effects

The full conditionals of the input effects (Θ, η) can be derived from the joint posterior (40) and are illustrated as follows.

(a) Conditioning on θ_2 and V_1 , the posterior distribution of θ_1 is

$$(\boldsymbol{\theta}_1 | \boldsymbol{\theta}_2, \boldsymbol{V}_1; \boldsymbol{y}) \sim N_{m_1 k_1} \left(\tilde{\boldsymbol{\theta}}_1, \tilde{\boldsymbol{V}}_1 \right),$$
 (41)

where

$$egin{aligned} & ilde{V}_1 = \left(oldsymbol{Z}_0^ op oldsymbol{\Sigma}^{-1} oldsymbol{Z}_0 + oldsymbol{I}_{m_1} \otimes oldsymbol{V}_1^{-1}
ight)^{-1}, \ & ilde{oldsymbol{ heta}}_1 = ilde{V}_1 igg[oldsymbol{Z}_0^ op oldsymbol{\Sigma}^{-1} oldsymbol{y} + \left(oldsymbol{I}_{m_1} \otimes oldsymbol{V}_1^{-1}
ight) oldsymbol{Z}_1 oldsymbol{ heta}_2 igg]. \end{aligned}$$

(b) The full conditional posteriors of θ_j, j = 2,..., r have the forms:

$$(\boldsymbol{\theta}_{j} | \boldsymbol{\theta}_{j-1}, \boldsymbol{\theta}_{j+1}, \boldsymbol{V}_{j-1}, \boldsymbol{V}_{j}) \sim N_{m_{j}k_{j}}\left(\tilde{\boldsymbol{\theta}}_{j}, \tilde{\boldsymbol{V}}_{j}\right),$$
(42)

where

$$\begin{split} \tilde{\boldsymbol{V}}_{j} &= \left[\boldsymbol{Z}_{j-1}^{\top} \left(\boldsymbol{I}_{m_{j-1}} \otimes \boldsymbol{V}_{j-1}^{-1} \right) \boldsymbol{Z}_{j-1} \right. \\ &+ \boldsymbol{I}_{m_{j}} \otimes \boldsymbol{V}_{j}^{-1} \right]^{-1}, \\ \tilde{\boldsymbol{\theta}}_{j} &= \tilde{\boldsymbol{V}}_{j} \left[\boldsymbol{Z}_{j-1}^{\top} \left(\boldsymbol{I}_{m_{j-1}} \otimes \boldsymbol{V}_{j-1}^{-1} \right) \boldsymbol{\theta}_{j-1} \right. \\ &+ \left(\boldsymbol{I}_{m_{j}} \otimes \boldsymbol{V}_{j}^{-1} \right) \boldsymbol{Z}_{j} \boldsymbol{\theta}_{j+1} \right], \end{split}$$

and $\boldsymbol{\theta}_{r+1} \triangleq \boldsymbol{\eta}$.

(c) By using (14), the full conditional of η can be derived as

$$(\boldsymbol{\eta} \mid \boldsymbol{\theta}_r, \lambda, \boldsymbol{V}_r) \sim N_d \left(\tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{V}}_{\boldsymbol{\eta}} \right),$$
 (43)

where

$$egin{aligned} & ilde{m{V}}_{m{\eta}} = \left[m{Z}_r^{ op} \left(m{I}_{m_r} \otimes m{V}_r^{-1}
ight) m{Z}_r + \lambda^{-1} m{I}_d
ight]^{-1}, \ & ilde{m{\eta}} = ilde{m{V}}_{m{\eta}} m{Z}_r^{ op} \left(m{I}_{m_r} \otimes m{V}_r^{-1}
ight) m{ heta}_r. \end{aligned}$$

Input effects $\theta_j \in \Theta$ and η can be readily sampled from their conditionals during the Gibbs sampling procedure, as their full conditional posterior distributions are all standard distributions.

The variance components which include V_j 's and λ can be updated from their full conditionals, and these conditionals have densities as follows.

(a) Given η , the conditional posterior density of λ is

$$\pi(\lambda|\boldsymbol{\eta}) \propto \lambda^{-\frac{d+1}{2}} \exp\left\{-\frac{1+\|\boldsymbol{\eta}\|^2}{2\lambda}\right\}, \quad (44)$$

which is actually an inverse gamma distribution as $IG\left(\frac{d-1}{2}, \frac{1+\|\eta\|^2}{2}\right)$.

(b) For $j \in [r]$, define that $t_j = \frac{m_j}{2} + 1 - \frac{1}{2k_j}$. The marginal posterior density of V_j given (θ_j, θ_{j+1}) is

$$\pi \left(\mathbf{V}_{j} | \boldsymbol{\theta}_{j}, \boldsymbol{\theta}_{j+1} \right) \propto \frac{1}{\left| \mathbf{V}_{j} \right|^{t_{j}} \prod_{1 \leq s < t \leq k_{j}} \left(\omega_{js} - \omega_{jt} \right)} \\ \times \operatorname{etr} \left\{ -\frac{1}{2} \mathbf{V}_{j}^{-1} \mathbf{H}_{j} \right\}, \qquad (45)$$

where etr(A) denotes exp(tr(A)) for a square matrix A, and

$$H_{j} \triangleq H_{j}(\boldsymbol{\theta}_{j}, \boldsymbol{\theta}_{j+1})$$

= $\sum_{i=1}^{m_{j}} (\boldsymbol{\theta}_{ji} - \boldsymbol{Z}_{ji} \boldsymbol{\theta}_{j+1}) (\boldsymbol{\theta}_{ji} - \boldsymbol{Z}_{ji} \boldsymbol{\theta}_{j+1})^{\top}$

The updating of λ can be simply carried out by sampling from an inverse gamma distribution. The full conditional posteriors of V_j , (45) are actually distributed as a recently proposed class of prior distributions by Berger et al. (2020a) for the covariance matrix, which is called the Shrinkage Inverse Wishart (SIW) distributions. The new class SIW(a, H) for a $k \times k$ covariance matrix W has the density as

$$\pi^{\text{SIW}}(\boldsymbol{W} \mid \boldsymbol{a}, \boldsymbol{H}) \propto \frac{\text{etr}\left(-\frac{1}{2}\boldsymbol{W}^{-1}\boldsymbol{H}\right)}{|\boldsymbol{W}|^{a}\prod_{i < j}\left(\nu_{i} - \nu_{j}\right)}, \qquad (46)$$

where $v_1 > v_2 > \cdots > v_k > 0$ are the ordered eigenvalues of W, a is a real constant and H is a $k \times k$ non-negative definite matrix. Thus, V_i are distributed as SIW $(t_i, H_i), j \in [r]$. To sample the covariance matrices from the full conditional posteriors, the previously suggested methods include the Metropolis-Hastings algorithm (cf. Berger et al., 2005) and Hit-and-run method (cf. Yang & Berger, 1994). The two methods both generate full candidate matrices by utilizing fullparameter proposal distributions, resulting in that they only work for moderate dimensions of the covariance matrices. To tackle this issue, Berger et al. (2020a) proposed a powerful Gibbs method for efficiently sampling the covariance matrices from their conditional densities and this new method works for higher dimensions k. The audience can refer to Berger et al. (2020a) or Appendix 3 for details of this Gibbs sampling

method. According to the simulation results of Berger et al. (2020a), the new Gibbs method outperforms the Metropolis-Hastings and Hit-and-run methods for moderate dimensions and work for k up to 100, while the other two algorithms break down in much lower dimensions.

In the framework of 2-level HNLM, Berger et al. (2020b) compared the numerical performance, from the mean square error (MSE) perspective, of a dozen of objective hyperpriors, which are the product of three objective hyperpriors for the hypermean and four objective hyperpriors for the hypercovariance matrix. Priors on the hypermean include constant prior, conjugate prior and the recommended prior (12). Priors on the hypercovariance matrix include constant prior, hierarchical Jefferys prior, hierarchical reference prior and the recommended prior (13). Their simulation results have shown that the recommended combination of hyperpriors dominates all the others in terms of Bayes risk, and the constant prior on the hypercovariance performs the worst. However, neither of the two remaining choices for hypercovariance is computationally easy. Considering the 4-level HNLM, Song et al. (2020) performed numerical experiment to compare the recommended prior with constant prior for the hypercovariance matrices, and the other two priors were canceled due to intractable computation. Also, Song et al. (2020)'s result presented the domination of the recommended hyperpriors over other priors. In conclusion, both Berger et al. (2020b) and Song et al. (2020) have provided strong numerical evidence of the superiority of the recommended hyperpriors for use in GHNLM, since 2-level and 4-level HNLMs both are specific GHNLM.

5. Discussions

We have proposed a generalized hierarchical normal linear model applicable to the nested data with complex structures. The GHNL model proves to be equivalent to a LMM model, while the GHNL model is more natural for researchers to model nested data from scratch, especially when incorporating covariates at high levels. Like generalizations to the simple normal linear model, the GHNL can be generalized to the hierarchical model with generalized liner model in the first level, and thus discrete observations can be handled. Besides, the first level (or even higher levels) of the GHNL model can be also extended to the setting of semiparametric regression models, such as, single index model and partially linear model. The technique of modelling and investigation in this paper can be applied to the linear part of the models mentioned above. The statistical analysis could be complicated and such explorations are beyond the scope of this paper, however.

Berger et al. (2020b) put an end to the endless search for the appropriate hyperpriors in hierarchical modelling and investigated properties comprehensively to justify the recommendation. Nonetheless, when it came to the propriety of the resulting posterior, they only suspected that that is true for use at any level for a general hierarchical normal model, the conditions for which were not given. To complete the story, we have studied the conditions for the posterior to be proper in more general situations than Berger et al. (2020b), when employing the recommended prior for the GHNL model. Theorems 3.1 and 3.3 demonstrate the main result, and Corollaries 3.1 and 3.3 reduce the computational burdens by defining kernel sets for $d \ge 2$ and d = 1, respectively. In addition, Theorems 3.2 and 3.4 provide powerful tools of simple forms for checking propriety of posterior for $d \ge 2$ and d = 1, separately. The user-friendly guidance for checking posterior propriety is eventually supplied. Note that our results only present sufficient conditions, and necessary conditions have never been discussed. The reason is because the derivation of the lower bound on the integrated likelihood of hyperparameters is intractable. Moreover, it is not worthwhile to investigate necessary conditions, as the derived upper bounds are tight enough such that the corresponding sufficient conditions are very modest, according to the remarks and examples in Section 3. At last, an efficient and powerful Gibbs sampling method for sampling from the posterior is introduced, overcoming the bottleneck of computation that the previously proposed sampling method only works for low dimensions or moderate dimensions inefficiently. The numerical evidence supporting the superiority of the recommended prior for hierarchical models was presented in Berger et al. (2020b) and Song et al. (2020).

Though we have made much progress in the hierarchical linear modelling, a major obstacle to applying our results is that the variance component for the first level is supposed to be known, which can hardly be satisfactory in practice. If we assume an unknown covariance matrix Σ_0 for the first level and specify prior (13) on it, the exponential term within the likelihood can not be dropped simply any longer when deriving the upper bound, otherwise, the resulting integral will be always infinite. The upper bound on the exponential term with respect to the eigenvalues of covariance matrices is very tricky to be obtained, and the condition for the integrability of the resulting integral remains to be further studied. Thus, the GHNL modelling with unknown Σ_0 can be taken as a sequential study of this paper.

Disclosure statement

No potential conflict of interest was reported by the author(s).

Funding

The research was supported by the National Natural Science Foundation of China [grant number 11671146].

References

- Berger, J. (1980). A robust generalized Bayes estimator and confidence region for a multivariate normal mean. *Annals* of Statistics, 8(4), 716–761. https://doi.org/10.1214/aos/ 1176345068
- Berger, J., Strawderman, W., & Tang, D. (2005). Posterior propriety and admissibility of hyperpriors in normal hierarchical models. *Annals of Statistics*, 33(2), 606–646. https://doi.org/10.1214/009053605000000075
- Berger, J., Sun, D., & Song, C. (2020a). Bayesian analysis of the covariance matrix of a multivariate normal distribution with a new class of priors. *Annals of Statistics*, 48(4), 2381–2403. https://doi.org/10.1214/19-AOS1891
- Berger, J., Sun, D., & Song, C. (2020b). An objective prior for hyperparameters in normal hierarchical models. *Journal of Multivariate Analysis*, 178(2020), 1–13. https://doi.org/10.1016/j.jmva.2020.104606
- Consonni, G., Fouskakis, D., Liseo, B., & Ntzoufras, I. (2018). Prior distributions for objective Bayesian analysis. *Bayesian Analysis*, 13(2), 627–679. https://doi.org/10.1214/ 18-BA1103
- Daniels, M. J., & Kass, R. E. (1999). Nonconjugate Bayesian estimation of covariance matrices and its use in hierarchical models. *Journal of the American Statistical Association.*, 94(448), 1254–1263. https://doi.org/10.1080/01621459.19 99.10473878
- Everson, P. J., & Morris, C. N. (2000). Inference for multivariate normal hierarchical models. *Journal of the Royal Statistical Society: Series B*, 62(2), 399–412. https://doi.org/10.11 11/rssb.2000.62.issue-2
- Fourdrinier, D., Strawderman, W. E., & Wells, M. T. (1998). On the construction of Bayes minimax estimators. *Annals of Statistics*, 26(2), 660–671. https://doi.org/10.1214/aos/ 1028144853
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, 1(3), 515–534. https://doi.org/10.1214/06-BA117A
- Goldstein, H. (2011). *Multilevel statistical models* (Vol. 922). John Wiley & Sons.
- Gustafson, P., Hossain, S., & Macnab, Y. C. (2006). Conservative prior distributions for variance parameters in hierarchical models. *Canadian Journal of Statistics*, 34(3), 377–390. https://doi.org/10.1002/cjs.v34:3
- Hobert, J. P., & Casella, G. (1996). The effect of improper priors on Gibbs sampling in hierarchical linear mixed models. *Journal of the American Statistical Association*, 91(436), 1461–1473. https://doi.org/10.1080/01621459.1996.1047 6714
- Hoff, P. D. (2009b). Simulation of the matrix Bingham-Von Mises-Fisher distribution, with applications to multivariate and relational data. *Journal of Computational* and Graphical Statistics, 18(2), 438–456. https://doi.org/ 10.1198/jcgs.2009.07177
- Horn, R. A., & Johnson, C. R. (2012). *Matrix analysis*. Cambridge university press.
- Lindenberger, U., & Pötter, U. (1998). The complex nature of unique and shared effects in hierarchical linear regression: Implications for developmental psychology. *Psychological Methods*, 3(2), 218–230. https://doi.org/10.1037/1082-989 X.3.2.218
- Michalak, S. E., & Morris, C. N. (2016). Posterior propriety for hierarchical models with log-Likelihoods that have norm bounds. *Bayesian Analysis*, 11(2), 545–571. https:// doi.org/10.1214/15-BA962
- Raudenbush, S., & Bryk, A. S. (1986). A hierarchical model for studying school effects. *Sociology of Education*, 59(1), 1–17. https://doi.org/10.2307/2112482

- Shimotsu, K. (2010). Exact local Whittle estimation of fractional integration with unknown mean and time trend. *Econometric Theory*, 26(2), 501–540. https://doi.org/ 10.1017/S0266466609100075
- Song, C., Sun, D., Fan, K., & Mu, R. (2020). Posterior propriety of an objective prior in a 4-Level normal hierarchical model. *Mathematical Problems in Engineering*, 2020. https://doi.org/10.1155/2020/8236934
- Speckman, P. L., & Sun, D. (2003). Fully Bayesian spline smoothing and intrinsic autoregressive priors. *Biometrika*, 90(2), 289–302. https://doi.org/10.1093/biomet/90.2.289
- Sun, D., Tsutakawa, R. K., & He, Z. (2001). Propriety of posteriors with improper priors in hierarchical linear mixed models. *Statistica Sinica*, 11(1), 77–95. http://www.jstor.org/stable/24306811
- Xia, A., Ma, H., & Carlin, B. P. (2011). Bayesian hierarchical modeling for detecting safety signals in clinical trials. *Journal of Biopharmaceutical Statistics*, 21(5), 1006–1029. https://doi.org/10.1080/10543406.2010.520181
- Yang, R., & Berger, J. (1994). Estimation of a covariance matrix using the reference prior. *Annals of Statistics*, 22(3), 1195–1211. https://doi.org/10.1214/aos/117632 5625

Appendices

Appendix 1. A special LMM

Consider a special LMM of the form:

$$y = X\beta + Z_1u_1 + \cdots + Z_ru_r + \epsilon, \qquad (A1)$$

where y denotes the observations and is an $n \times 1$ vector, β is the vector of fixed effects and is an $p \times 1$ vector. For $i \in [r]$, u_i 's are $q_i \times 1$ vectors and represent the vectors of random effects, and u_i 's are assumed to be independently distributed as $N_{q_i}(\mathbf{0}, W_i)$, where W_i 's are $q_i \times q_i$ positive definite matrices and unknown. X is an $n \times p$ matrix, Z_i 's are $n \times q_i$ matrices, and X and Z_i 's are known design matrices. ϵ is the vector of random errors and distributed as $N_n(\mathbf{0}, \Sigma)$, Σ is an $n \times n$ positive definite matrix and given.

It follows from Berger et al. (2020b) that we can assume independent priors on $(\beta, W_1, \ldots, W_r)$ as

$$\pi(\boldsymbol{\beta}) \propto \frac{1}{(1+\|\boldsymbol{\beta}\|^2)^{(p-1)/2}}, \quad \boldsymbol{\beta} \in \mathbb{R}^p,$$

$$\pi(W_j) \propto \frac{1}{|W_j|^{1-1/(2q_j)} \prod_{1 \le s < t \le q_j} (\nu_{js} - \nu_{jt})},$$

$$W_j > 0, \ j \in [r], \tag{A2}$$

where $v_{j1} > v_{j2} > \cdots > v_{jq_j} > 0$ are the ordered eigenvalues of $W_j, j \in [r]$. The prior on β has a hierarchical structure of the form

$$(\boldsymbol{\beta}|\tau) \sim N_d(\boldsymbol{0},\tau \boldsymbol{I}_d) \text{ and } [\tau] \propto \tau^{-1/2} \exp\left(-\frac{1}{2\tau}\right).$$

The posterior propriety results for the special LMM (A1) is displayed as follows. Firstly, let $\tau \triangleq v_{01}$ and $q_0 = 1$. Denote the index set of the variance scale or the eigenvalues of the covariance matrices by $F = \{(j, l) | j = 0, 1, ..., r, l \in [q_j]\}$, and $\mathcal{T} = \{D | D \subseteq F, D \neq \emptyset\}$ represents the set of the nonempty subsets of *F*. Define that $c_{jl,s} = 1_{\{l=s\}}$ for $j \in [r], l \in [q_j]$ and $s \in [n]$.

Theorem A.1: Consider linear mixed effect model (A1) with prior (A2) on $(\beta, W_1, ..., W_r)$. Assume p > 1, and then the

posterior is proper if

$$\sum_{s=1}^{n} 1_{\left\{1_{\{(0,1)\in D\}}1_{\{s\leq p\}}+\sum_{j\neq 0, (j,l)\in D} c_{jl,s}>0\right\}} > \sum_{(j,l)\in D} \frac{1}{q_{j}}$$
(A3)

holds for any $D \in T$.

The proof of Theorem A.1 is similar to that of Theorem 3.1 and is omitted here.

Fact A.1: When p > 1, (A3) holds for any $D \in T$ iff

$$\sum_{j \in [r]} \frac{1}{q_j} < 1 \quad \text{and} \quad p > 1 + \sum_{j \in [r]} \frac{\min(p, q_j)}{q_j}.$$
 (A4)

Proof: For any $D \in \mathcal{T}$ and $(0, 1) \notin D$, (A3) is equivalent to

$$\sum_{s=1}^{n} 1_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > \sum_{(j,l)\in D} \frac{1}{q_j}.$$
 (A5)

It can be deduced that inequality (A5) holds for any $D \in \mathcal{T}$ and $(0, 1) \notin D$ iff

$$L > \sum_{j \in [r]} \frac{\min(L, q_j)}{q_j}, \quad \text{for } L \in [n],$$
 (A6)

which is equivalent to $\sum_{j \in [r]} \frac{1}{q_j} < 1$ since $q_j \ge 1$ for $j \in [r]$. Inequality (A3) holds for any $D \in \mathcal{T}$ with $(0, 1) \in D$ iff

$$L > 1 + \sum_{j \in [r]} \frac{\min(L, q_j)}{q_j}, \quad \text{for } L = p, \dots, n.$$
 (A7)

Under the condition that $\sum_{j \in [r]} \frac{1}{q_j} < 1$, (A7) is equivalent to

$$p>1+\sum_{j\in[r]}\frac{\min(p,q_j)}{q_j}.$$

Corollary A.1: Consider model (A1) with prior (A2) on parameters. The posterior is proper if one of the following condition holds,

(a)
$$p > l + r$$
 and $\sum_{j \in [r]} \frac{1}{q_j} < 1;$
(b) $p > l$ and $\sum_{j \in [r]} \frac{1}{q_j} < 1 - \frac{1}{p}.$

Proof: Since

$$\sum_{j\in[r]}\frac{\min(p,q_j)}{q_j}\leq r \quad \text{and} \quad \sum_{j\in[r]}\frac{\min(p,q_j)}{q_j}\leq p\sum_{j\in[r]}\frac{1}{q_j},$$

(a) and (b) follow from Fact A.1 directly.

Remark A.1: Consider model (A1) with r = 1. The posterior using prior (A2) is prior if either (a) $p \ge 2$, $q_1 \ge 3$ or (b) $p \ge 3$, $q_1 \ge 2$ holds. If p = 1, the posterior propriety can hardly be satisfied, the reason of which is two-fold. First, inequality (A3) fails for $D = \{(0, 1)\}$. Second, if we follow the thread of deriving the condition in Theorem 3.3, a sufficient condition can be derived as

$$\sum_{s=1}^{n} \mathbb{1}_{\left\{\sum_{(j,l)\in D} c_{jl,s} > 0\right\}} > \mathbb{1}_{\left\{\exists j \in [r], (j,1)\in D\right\}} + \sum_{(j,l)\in D} \frac{1}{q_j}, \quad (A8)$$

for any $D \in \mathcal{T}$ with $(0, 1) \notin D$. However, inequality (A8) does not hold for $D = \{(j, 1)\}, j \in [r]$.

Appendix 2. Proof of lemmas in Section 3.1

Lemma A.1: *min-max theorem, cf. Horn & Johnson,* (2012): For an $n \times n$ symmetric matrix A and a non-zero $n \times 1$ vector x, the Rayleigh quotient for A and x can be defined as

$$R(A, \mathbf{x}) = \frac{\langle A\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle},$$

where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product. Then

$$\lambda_k(A) = \max_U \left\{ \min_{\boldsymbol{x}} \left\{ R(A, \boldsymbol{x}) \mid \boldsymbol{x} \in U, \\ \|\boldsymbol{x}\| \neq 0 \right\} \mid \dim(U) = k \right\}, \quad k \in [n]$$
(A9)

where U denotes the linear subspace of \mathbb{R}^n . Especially,

$$\lambda_{\max}(\mathbf{A}) = \max_{\mathbf{x}} \{ R(\mathbf{A}, \mathbf{x}) \mid \|\mathbf{x}\| \neq 0 \} \text{ and}$$
$$\lambda_{\min} = \min_{\mathbf{x}}(\mathbf{A}) \{ R(\mathbf{A}, \mathbf{x}) \mid \|\mathbf{x}\| \neq 0 \}.$$
(A10)

Lemma A.2: For $n \times n$ symmetric matrices A_j , $j \in [r]$, we have

(a) $\lambda_{\max}\left(\sum_{j=1}^{r} A_{j}\right) \leq \sum_{j=1}^{r} \lambda_{\max}(A_{j});$ (b) supposing $A_{j} \geq 0, j \in [r], then$

$$\lambda_k\left(\sum_{j=1}^r A_j\right) \ge \frac{1}{r} \sum_{j=1}^r \lambda_k\left(A_j\right), \quad k \in [n]$$

Proof: For (a), given an $n \times 1$ non-zero vector x, by (A10),

$$R\left(\sum_{j=1}^{r} A_{j}, \mathbf{x}\right) = \sum_{j=1}^{r} R(A_{j}, \mathbf{x}) \leq \sum_{j=1}^{r} \max_{\mathbf{x}_{j} \neq \mathbf{0}} R(A_{j}, \mathbf{x}_{j})$$
$$= \sum_{j=1}^{r} \lambda_{\max}(A_{j}).$$
(A11)

The proof for (a) is completed by using (A10) again. For (b), it suffices to prove that for any $j \in [r]$

$$\lambda_k\left(\sum_{j=1}^r A_j\right) \ge \lambda_k(A_j), \quad i \in [n].$$
 (A12)

Since $A_l \ge 0$, $l \in [r]$, for any $j \in [r]$, $x \in \mathbb{R}^n$ and $x \neq 0$,

$$R\left(\sum_{l=1}^{r} A_{l}, \boldsymbol{x}\right) = \sum_{l=1}^{r} R(A_{l}, \boldsymbol{x}) \ge R(A_{j}, \boldsymbol{x}).$$

Minimize the both sides of the inequality above over $\{x \mid x \in U, x \neq 0\}$ first. Then take the maximum over $\{U \mid U \subseteq \mathbb{R}^n, \dim(U) = k\}$. (A12) can be easily obtained by using Lemma A.1.

A.1 Proof of Lemma 3.1

For (a), by Lemma A.2 (a), it suffices to prove that $\lambda_{\max}(\mathbf{X}_j \mathbf{A}_j \mathbf{X}_j^{\top}) \leq C_1 \lambda_{\max}(\mathbf{A}_j), j \in [r]$. For any $j \in [r]$, applying (A10), yields

$$0 < \lambda_{\max}(X_j A_j X_j^{\top}) = R\left(X_j A_j X_j^{\top}, x^*\right)$$
$$= \max_{x \neq 0} R\left(X_j A_j X_j^{\top}, x\right).$$

It is obvious that $X_j^{\top} \mathbf{x}^* \neq \mathbf{0}$, otherwise, it will result in $R\left(X_j A_j X_j^{\top}, \mathbf{x}^*\right) = 0$, which contradicts. In addition, since

 $\left\langle \boldsymbol{X}_{j}^{\top}\boldsymbol{x}^{*}, \boldsymbol{X}_{j}^{\top}\boldsymbol{x}^{*} \right\rangle \leq \lambda_{\max}\left(\boldsymbol{X}_{j}^{\top}\boldsymbol{X}_{j}\right) \left\langle \boldsymbol{x}^{*}, \boldsymbol{x}^{*} \right\rangle \leq C_{1}\left\langle \boldsymbol{x}^{*}, \boldsymbol{x}^{*} \right\rangle, \quad \text{we}$ have

$$R\left(\boldsymbol{X}_{j}\boldsymbol{A}_{j}\boldsymbol{X}_{j}^{\top},\boldsymbol{x}^{*}\right) \leq C_{1}R\left(\boldsymbol{A}_{j},\boldsymbol{X}_{j}^{\top}\boldsymbol{x}^{*}\right)$$
$$\leq C_{1}\max\left\{R\left(\boldsymbol{A}_{j},\boldsymbol{z}\right) \mid \boldsymbol{z} \in \mathbb{R}^{p_{j}}, \, \boldsymbol{z} \neq \boldsymbol{0}\right\}$$
$$= C_{1}\lambda_{\max}(\boldsymbol{A}_{j}).$$

Therefore, we have proved part (a).

2

For (b), we only need to prove that

$$C_k(\boldsymbol{H}) \geq \frac{C_2}{r} \sum_{j=1}^r a_{jk}, \quad k \in [n].$$

It follows from Lemma A.2 (b) that for any $k \in [n]$

$$\lambda_k(\boldsymbol{H}) \geq \frac{1}{r} \sum_{j=1}^r \lambda_k(\boldsymbol{X}_j \boldsymbol{A}_j \boldsymbol{X}_j^{\top}).$$

Since rank $(X_j A_j X_j^{\top})$ = rank (A_j) , then $\lambda_k (X_j A_j X_j^{\top}) = 0, p_j < k \le n$. Thus, it remains to show that

$$\lambda_k(\mathbf{X}_j \mathbf{A}_j \mathbf{X}_j^{\top}) \ge C_2 \lambda_k(\mathbf{A}_j), \quad j \in [r], \ k \in [p_j].$$
(A13)

Firstly, for any $j \in [r]$, we introduce a linear transformation $\mathcal{L}_j : \mathbb{R}^n \to \mathbb{R}^{p_j}$ defined as $\mathcal{L}_j(\mathbf{v}) = \mathbf{X}_j^\top \mathbf{v}, \mathbf{v} \in \mathbb{R}^n$. The *kernel* space of \mathcal{L}_j is denoted by $\operatorname{Ker}(\mathcal{L}_j) = \{\mathbf{v} \in \mathbb{R}^n : \mathcal{L}_j(\mathbf{v}) = \mathbf{0}\}$. Since \mathbf{X}_j is of full column rank $p_j \leq n$, then the dimension of the complementary space of $\operatorname{Ker}(\mathcal{L}_j)$ is p_j , i.e., $\dim (\operatorname{Ker}(\mathcal{L}_j)^{\perp}) = p_j$. Thus, the mapping $\mathcal{L}_j : \operatorname{Ker}(\mathcal{L}_j)^{\perp} \mapsto \mathbb{R}^{p_j}$ is a one-to-one mapping. For any $U \subseteq \mathbb{R}^{p_j}$ with $\dim(U) = k, k \in [p_j]$, define that

$$\mathcal{L}_j^*(U) = \left\{ \boldsymbol{v} \in \operatorname{Ker}(\mathcal{L}_j)^{\perp} : \mathcal{L}_j(\boldsymbol{v}) = \boldsymbol{x}, \quad \boldsymbol{x} \in U \right\}.$$

It is obvious that $\mathcal{L}_{j}^{*}(U) \subseteq \operatorname{Ker}(\mathcal{L}_{j})^{\perp}$ and dim $(\mathcal{L}_{j}^{*}(U)) = k$. For any $U \subseteq \mathbb{R}^{p_{j}}$ with dim(U) = k and any $\mathbf{x} \in U$, there exists one and only one $\mathbf{v} \in \mathcal{L}_{j}^{*}(U)$ such that $\mathcal{L}_{j}(v) = \mathbf{x}$. Since $\langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{v}^{\top}(\mathbf{X}_{j}^{\top}\mathbf{X}_{j})\mathbf{v} \geq C_{2} \langle \mathbf{v}, \mathbf{v} \rangle$, we have

$$R\left(\boldsymbol{X}_{j}\boldsymbol{A}_{j}\boldsymbol{X}_{j}^{\top}, \boldsymbol{\nu}\right) \geq C_{2}R\left(\boldsymbol{A}_{j}, \boldsymbol{x}\right).$$
(A14)

It refers to Lemma A.1 that

$$\lambda_{k}(\boldsymbol{X}_{j}\boldsymbol{A}_{j}\boldsymbol{X}_{j}^{\top}) = \max_{V} \left\{ \min_{\boldsymbol{\nu}} \left\{ R(\boldsymbol{X}_{j}\boldsymbol{A}_{j}\boldsymbol{X}_{j}^{\top},\boldsymbol{\nu}) \mid \boldsymbol{\nu} \in V, \|\boldsymbol{\nu}\| \neq 0 \right\} \\ \left| V \subseteq \operatorname{Ker}(\mathcal{L}_{j})^{\perp}, \operatorname{dim}(V) = k \right\},$$
(A15)

for $k \in [p_j]$. Minimize the both sides of inequality (A14) over $\{\boldsymbol{x} \mid \boldsymbol{x} \in U, \boldsymbol{x} \neq \boldsymbol{0}\}$ first. Then take the maximum over $\{U \mid U \subseteq \mathbb{R}^{p_j}, \dim(U) = k\}$, (A13) can be easily obtained by using Lemma A.1 and (A15).

A.2 Proof of Lemma 3.2

The Domain of the integral can be divided into

$$\Omega_0 = \{ \boldsymbol{\lambda} \mid 0 \le \lambda_j \le 1, \ j \in [k] \},$$

$$\Omega_D = \{ \boldsymbol{\lambda} \mid \lambda_j > 1, \ j \in D \text{ and } 0 \le \lambda_i \le 1, \ i \in [k]/D \},$$

$$D \in \mathcal{F}([k])$$

i.e., $\Omega = \left(\bigcup_{D \in \mathcal{F}([k])} \Omega_D\right) \bigcup \Omega_0$. Thus, the integral *I* is finite iff the integrals over Ω_0 and Ω_D for each $D \in \mathcal{F}([k])$ are finite.

Denote the integrand as $F(\lambda)$. Then

$$\int_{\Omega_0} F(oldsymbol{\lambda}) \, \mathrm{d}oldsymbol{\lambda} \simeq \int_{\Omega_0} rac{1}{\prod_{j=1}^k \lambda_j^{a_j}} \, \mathrm{d}oldsymbol{\lambda},$$

which is finite iff condition (a) is satisfied.

To verify condition (b), we only need to justify the following statement. Also, we assume condition (a) is always satisfied hereafter.

Fact A.2: For all $D \in \mathcal{F}([k])$ with $\operatorname{card}(D) = L$, $1 \le L \le k$, the integrals $\int_{\Omega_D} F(\lambda) d\lambda$ are finite iff inequalities

$$\sum_{j\in E} a_j + \sum_{i\in \mathcal{G}_E} b_i > \operatorname{card}(E)$$
(A16)

hold for all $E \in \mathcal{F}([k])$ with card(E) $\leq L$ and $\mathcal{G}_E = \{i | E \cap C_i \neq \emptyset, i \in [n]\}$. Under the condition above,

$$\int_{\Theta_D(t)} G_D(\lambda) \left(\prod_{r \in D} d\lambda_r \right)$$
$$\simeq \exp\left\{ -\log t \left(\sum_{j \in D} a_j + \sum_{i \in \mathcal{G}_D} b_i - \operatorname{card}(D) \right) \right\}$$
(A17)

always holds, where

$$\Theta_D(t) = \left\{ \boldsymbol{\lambda} \mid \lambda_j \ge t, \ j \in D \quad \text{and} \quad 0 \le \lambda_i \le 1, \ i \in [k]/D \right\},$$
$$G_D(\boldsymbol{\lambda}) = \frac{1}{\left[\prod_{j \in D} \lambda_j^{a_j} \right] \left[\prod_{i \in \mathcal{G}_D} (1 + \sum_{r \in C_i} \lambda_r)^{b_i} \right]}.$$

The reason why formula (A17) is required is that it plays an important role in verifying condition (A16).

Proof: We prove the result by the technique of mathematical induction. First, we assume that the statement in Fact A.2 is true for L = l, $1 \le l \le (k - 1)$. With this assumption, we must show that the statement is true for its successor, L = (l + 1). Write an arbitrary set $D \in \mathcal{F}([k])$ with cardinality (l + 1) as $\{j_1, \ldots, j_{l+1}\}$, where $1 \le j_1 < \cdots < j_{l+1} \le k$. Denote that $D_{-j_i} = D/[\{j_i\}], i = 1, \ldots, (l + 1)$.

Step 1: We first prove that $\int_{\Omega_D} F(\lambda) d\lambda$ is finite iff inequalities (A16) hold for L = (l + 1).

Region Ω_D can be divided into

$$\Sigma_{1} = \left\{ \boldsymbol{\lambda} | \lambda_{j} \geq \lambda_{j_{1}} > 1, j \in D_{-j_{1}} \text{ and} \\ 0 \leq \lambda_{i} \leq 1, i \in [k]/D \right\}$$

$$\vdots$$

$$\Sigma_{l+1} = \left\{ \boldsymbol{\lambda} | \lambda_{j} \geq \lambda_{j_{l+1}} > 1, j \in D_{-j_{l+1}} \\ \text{and } 0 \leq \lambda_{i} \leq 1, i \in [k]/D \right\}.$$

Therefore, integral $\int_{\Omega_D} F(\lambda) \, d\lambda$ is finite iff $\int_{\Sigma_i} F(\lambda) \, d\lambda < \infty$ for any $i = 1, \dots, (l+1)$. For $i = 1, \dots, (l+1)$, we have

$$\int_{\Sigma_i} F(\boldsymbol{\lambda}) \, \mathrm{d}\boldsymbol{\lambda} \simeq \int_{\Sigma_i} \frac{1}{\prod_{s \notin D} \lambda_s^{a_s}} \left(\frac{1}{\lambda_{j_i}} \right)^{a_{j_i} + \sum_{r \in \mathcal{H}_i} b_r} G_{D_{-j_i}}(\boldsymbol{\lambda}) \, \mathrm{d}\boldsymbol{\lambda},$$

where $\mathcal{H}_i = \{r \mid C_r \cap D_{j_i} = \emptyset \text{ and } j_i \in C_r, r \in [n]\}$, and it's easy to see that $\mathcal{G}_D = \mathcal{G}_{D_{-j_i}} \bigcup \mathcal{H}_i$ and $\mathcal{G}_{D_{-j_i}} \cap \mathcal{H}_i = \emptyset$. Since

$$\int_{\Sigma_i} G_{D_{-j_i}}(\boldsymbol{\lambda}) \left(\prod_{r \in D_{-j_i}} \mathrm{d} \lambda_r \right)$$

$$\simeq \int_{\Theta_{D_{-j_i}}(\lambda_{j_i})} G_{D_{-j_i}}(\lambda) \left(\prod_{r \in D_{-j_i}} d\lambda_r\right), \qquad (A18)$$

$$\int_{\Omega_{D-j_i}} F(\boldsymbol{\lambda}) \, \mathrm{d}\boldsymbol{\lambda} \simeq \int_{\Omega_{D-j_i}} \frac{1}{\prod_{s \notin D-j_i} \lambda_s^{a_s}} G_{D-j_i}(\boldsymbol{\lambda}) \, \mathrm{d}\boldsymbol{\lambda} \quad (A19)$$

and the RHS (Right-Hand Side) of (A19) and (A18) are finite simultaneously under condition (a). Hence, The LHS (Left-Hand Side) of (A18) is finite iff $\int_{\Omega_{D_{-j_i}}} F(\lambda) d\lambda$ is finite.

Furthermore, by assumption and (A17), we have

$$\begin{split} &\int_{\Sigma_i} G_{D_{-j_i}}(\lambda) \left(\prod_{r \in D_{-j_i}} d\lambda_r\right) \\ &\simeq \exp\left\{-\log \lambda_{j_i} \left(\sum_{j \in D_{-j_i}} a_j + \sum_{r \in \mathcal{G}_{D_{-j_i}}} b_r - l\right)\right\}. \end{split}$$

Thus, under condition (a) and assumption, we have

$$\int_{\Sigma_i} F(\boldsymbol{\lambda}) \, \mathrm{d}\boldsymbol{\lambda} \simeq \int_1^\infty \left(\frac{1}{\lambda_{j_i}}\right)^{\sum_{j \in D} a_j + \sum_{r \in \mathcal{G}_D} b_r - l} \, \mathrm{d}\lambda_{j_i},$$

the RHS of which is finite iff $\sum_{j \in D} a_j + \sum_{r \in \mathcal{G}_D} b_r > 1 + l = card(D)$.

In conclusion, $\int_{\Omega_D} F(\lambda) d\lambda$ is finite iff $\sum_{j \in D} a_j + \sum_{r \in \mathcal{G}_D} b_r > \operatorname{card}(D)$ and $\int_{\Omega_{D-j_i}} F(\lambda) d\lambda$ is finite for any $i \in [l+1]$. Since *D* is arbitrary and $\operatorname{card}(D_{-j_i}) = l$, we have accomplished the goal of Step 1.

Step 2: Next, we prove that formula (A17) holds for *D* with cardinality (l + 1).

Region $\Theta_D(t)$ can be divided into

(1)

$$\Theta_D^{(1)}(t) = \left\{ \boldsymbol{\lambda} | \lambda_j \ge \lambda_{j_1} \ge t, j \in D_{-j_1} \text{ and} \right.$$
$$0 \le \lambda_i \le 1, i \in [k]/D \right\}$$
$$\vdots$$
$$\Theta_D^{(l+1)}(t) = \left\{ \boldsymbol{\lambda} | \lambda_j \ge \lambda_{j_{l+1}} \ge t, j \in D_{-j_{l+1}} \text{ and} \right.$$

Similar to the proof of Step 1, we can prove that for i = 1, ..., (l + 1),

 $0 \leq \lambda_i \leq 1, i \in [k]/D\}.$

$$\begin{split} &\int_{\Theta_D^{(i)}(t)} G_D(\lambda) \left(\prod_{r \in D} d\lambda_r \right) \\ &\simeq \int_t^\infty \left(\frac{1}{\lambda_{j_i}} \right)^{\sum_{j \in D} a_j + \sum_{r \in \mathcal{G}_D} b_r - l} d\lambda_{j_i} \\ &= \exp\left\{ -\log t \left(\sum_{j \in D} a_j + \sum_{r \in \mathcal{G}_D} b_r - \operatorname{card}(D) \right) \right\}. \end{split}$$

Therefore, we get Step 2 proved.

Step 3: We need to present that the statement is true for L = 1 to complete the proof, on the basis of mathematical induction.

Denote that $D = \{r\}, r = 1, \dots, k$. Then

$$\int_{\Omega_D} F(\boldsymbol{\lambda}) \, \mathrm{d} \boldsymbol{\lambda} \simeq \int_1^\infty \left(\frac{1}{\lambda_r}\right)^{a_r + \sum_{i \in \mathcal{G}_D} b_i} \, \mathrm{d} \lambda_r$$

which is finite iff $\sum_{j \in D} a_j + \sum_{i \in \mathcal{G}_D} b_i > 1 = \operatorname{card}(D)$, under which,

$$\begin{split} &\int_{\Theta_D(t)} G_D(\lambda) \left(\prod_{i \in D} d\lambda_i \right) \\ &\simeq \int_t^\infty \left(\frac{1}{\lambda_r} \right)^{a_r + \sum_{i \in \mathcal{G}_D} b_i} d\lambda_r \\ &= \exp\left\{ -\log t \left(\sum_{j \in D} a_j + \sum_{i \in \mathcal{G}_D} b_i - \operatorname{card}(D) \right) \right\}, \end{split}$$

which accomplishes the proof of Fact A.2.

Appendix 3. Gibbs sampling from the SIW distributions

As for the SIW distribution (46), we first consider the change of variables from W to $\Xi = \text{diag}(v_1, \ldots, v_k)$ and the orthogonal matrix O of corresponding eigenvectors. The Jacobian is

$$\left|\frac{\partial W}{\partial (\boldsymbol{\Xi}, \boldsymbol{O})}\right| = \prod_{i < j} (v_i - v_j).$$
 (A20)

According to (A20) and Lemma 4 in Berger et al. (2020a), (46) can be transformed to

$$\pi(\Xi, \mathbf{O}) \propto \frac{1}{|\Xi|^a} \operatorname{etr} \left(-\frac{1}{2} \Xi^{-1} \mathbf{O}' H \mathbf{O} \right).$$
 (A21)

Gibbs sampling of Ξ : We first sample Ξ given (O, H) from

$$\pi(\Xi \mid \boldsymbol{O}, \boldsymbol{H}) \propto \frac{1}{\prod_{i=1}^{k} v_i^a} \operatorname{etr}\left(-\frac{1}{2}\Xi^{-1}\boldsymbol{O}'\boldsymbol{H}\boldsymbol{O}\right)$$
$$= \prod_{i=1}^{k} \frac{1}{v_i^a} \operatorname{etr}\left(-\frac{c_i}{v_i}\right),$$

where c_i is the *i*-th diagonal element of O'HO/2, $i \in [k]$. Therefore, we can sample v_i independently from $IG(a-1,c_i)$.

Gibbs sampling of O: Given (Ξ, H) , the marginal density of **O** has the form:

$$\pi(\mathbf{O} \mid \mathbf{\Xi}, \mathbf{H}) \propto \operatorname{etr}\left(-\frac{1}{2}\mathbf{H}\mathbf{O}\mathbf{\Xi}^{-1}\mathbf{O}'\right).$$

Let $H = LUL^{\top}$, where $LL^{\top} = I_k$ and $U = \text{diag}(u_1, \dots, u_k)$ is the diagonal matrix of corresponding eigenvalues with

 $u_1 \ge \cdots \ge u_k$. Define $G = L^{\top}O$. Since the invariant right Haar measure is invariant to the orthonormal transformation, the conditional density of *G* is

$$\pi(\boldsymbol{G} \mid \boldsymbol{\Xi}, \boldsymbol{H}) \propto \operatorname{etr}\left(-\frac{1}{2}\boldsymbol{U}\boldsymbol{G}\boldsymbol{\Xi}^{-1}\boldsymbol{G}'\right).$$
 (A22)

The updating of *G* from (A22) can be implemented by applying a Gibbs update to two randomly selected columns (cf. Hoff, 2009b) or rows (cf. Berger et al., 2020a). The two ways are essentially equivalent when rank(H) = k, but Berger et al. (2020a)'s method is considerably faster if rank(H) < k. Without any loss, assume that the two randomly selected rows are the first and second rows. The updated value of *G* can be written as $G^{\text{new}} = \text{diag} \left(\Phi, I_{k-2} \right) \begin{pmatrix} G_{12}^{\text{old}} \\ G_{-12}^{\text{old}} \end{pmatrix}$, where G_{12}^{old} denotes the first two rows of the old value of *G* which is G^{old} , G^{old}_{-12} is the remaining k-2 rows of G^{old} and

$$\mathbf{\Phi} = \mathbf{D}_{\epsilon} \mathbf{\Phi}_0 = \begin{pmatrix} \epsilon_1 & 0 \\ 0 & \epsilon_2 \end{pmatrix} \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix},$$

with $\phi \in (-\frac{\pi}{2}, \frac{\pi}{2}]$ and $\epsilon_i = \pm 1$ for i = 1, 2. Let $U_1 = \text{diag}(u_1, u_2)$. The full conditional density of ϕ has the form:

$$\pi(\phi \mid \boldsymbol{G}^{\text{old}}, \boldsymbol{\Xi}, \boldsymbol{H}) \propto \operatorname{etr} \left\{ -\frac{1}{2} \boldsymbol{U}_1 \boldsymbol{\Phi}_0 \boldsymbol{G}_{12}^{\text{old}} \boldsymbol{\Xi}^{-1} \left(\boldsymbol{G}_{12}^{\text{old}} \right)^\top \boldsymbol{\Phi}_0^\top \right\}.$$

Write

$$G_{12}^{\text{old}} \Xi^{-1} \left(G_{12}^{\text{old}} \right)^{\top} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} s_1 & 0 \\ 0 & s_2 \end{pmatrix} \times \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

where $\theta \in (-\frac{\pi}{2}, \frac{\pi}{2}]$ and $s_1 > s_2$. Then the conditional density of ϕ can be rewritten as

$$\pi(\phi \mid \mathbf{G}^{\text{old}}, \boldsymbol{\Xi}, \boldsymbol{H}) \propto \exp\left\{-c_0 \cos^2(\phi + \theta)\right\},$$

where $c_0 = \frac{1}{2}(s_1 - s_2)(u_1 - u_2) \ge 0$. Define $\alpha = \cos^2(\phi + \theta)$. Then the full conditional density of α has the form:

$$\pi(\alpha \mid \boldsymbol{G}^{\text{old}}, \boldsymbol{\Xi}, \boldsymbol{H}) \propto \exp\{-c_0 \alpha\} \alpha^{-\frac{1}{2}} (1-\alpha)^{-\frac{1}{2}}, \\ \alpha \in [0, 1].$$

Simulating $\alpha \in [0, 1]$ can proceed with a rejection sampler by setting the proposal distribution as Beta $(\frac{1}{2}, \frac{1}{2})$.