

# Joint Estimation and Inference in Data Integration with Multiple Multi-layered Gaussian Graphical Models

Subhabrata Majumdar, George Michailidis

**Abstract:** The rapid development of high-throughput technologies has enabled generation of data from biological processes that span multiple layers, like genomic, proteomic or metabolomic data; and pertain to multiple sources, like disease subtypes or experimental conditions. In this work we propose a general statistical framework based on graphical models for horizontal (i.e. across conditions or subtypes) and vertical (i.e. across different layers containing data on molecular compartments) integration of information in such datasets. We start with decomposing the multi-layer problem into a series of two-layer problems. For each two-layer problem, we model the outcomes at a node in the lower layer as dependent on those of other nodes in that layer, as well as all nodes in the upper layer. We use a combination of neighborhood selection and group-penalized regression to obtain sparse estimates of all model parameters. Following this, we propose a debiasing technique and asymptotic distributions of inter-layer directed edge weights that utilize already computed neighborhood selection coefficients for nodes in the upper layer. Subsequently we establish global and simultaneous testing procedures for these edge weights. Performance of our proposed methodology is analyzed using simulation experiments.

**Keywords:** Data integration; Gaussian Graphical Models; Neighborhood selection; Group lasso; hypothesis testing; multiple testing; false discovery rate

# 1 Introduction

The human body is a complex system, comprising of numerous regulatory networks that are connected within and between themselves. These networks have a natural hierarchy, for example DNA variations (Copy Number variations, Single Nucleotide Polymorphisms) influence gene expressions, which in turn influence RNA and protein expressions. They can also come from multiple sources that have some similarity with each other, for example diseased and healthy individuals, different experimental conditions, or different subtypes of the same disease. Recent developments in high-throughput technologies have resulted in the generation of enormous amount of data covering each of the above and many more situations (e.g. The Cancer Genome Atlas (TCGA: [Tomczak et al. \(2015\)](#))).

Figure 1 gives a schematic representation of the horizontal and vertical structure of heterogeneous biological data as outlined above. A simultaneous analysis of all parameters in this complex layered structure is known as *data integration*. While it is common knowledge that this will result in a more comprehensive picture of the regulatory mechanisms behind diseases, phenotypes and biological processes in general, there is a dearth of rigorous methodologies that satisfactorily tackle all challenges that stem from attempts to perform data integration ([Gligorićević and Pržulj, 2015](#); [Gomez-Cabrero et al., 2014](#); [Joyce and Palsson, 2006](#)). A review of the present approaches towards achieving this goal, which are based mostly on specific case studies, can be found in [Gligorićević and Pržulj \(2015\)](#) and [Zhang et al. \(2017\)](#).

Gaussian Graphical Models (GGM) have been extensively used to model biological networks in the last few years. While the initial work on GGMs focused on estimating undirected edges within a single network through obtaining sparse estimates of the inverse covariance matrix, or the precision matrix from the data (e.g. see the references in [Bühlmann and van de Geer \(2011\)](#)), the attention has now shifted to estimating parameters from more complex structures, especially multiple related graphical models and hierarchical multilayer models with both directed and undirected edges. For the first problem, [Guo et al. \(2011\)](#) and [Xie et al. \(2016\)](#) assumed perturbations over a common underlying structure to model multiple precision matrices, while [Danaher et al. \(2014\)](#) proposed using fused lasso or group lasso penalties in a joint group lasso model for the same purpose. To incorporate prior information on the group structures across several graphs, [Ma and Michailidis \(2016\)](#) proposed the Joint Structural Estimation Method (JSEM), which uses group-penalized neighborhood regression and subsequent refitting for estimating precision matrices.

For the second problem, a two-layered structure can be modeled by interpreting directed edges between the two layers as elements of a multitask regression coefficient matrix, while undirected edges inside either layer correspond to the precision matrix of predictors in that layer. While several methods exist in the literature for joint estimation of both parameters ([Cai et al., 2012](#); [Lee and Liu, 2012](#); [Rothman et al., 2010](#)), only recently [Lin et al. \(2016\)](#) made the observation that a multi-layer model can, in fact, be decomposed into a series of two-layer problems. Subsequently, they proposed an estimation algorithm and derived theoretical properties of the resulting estimators.

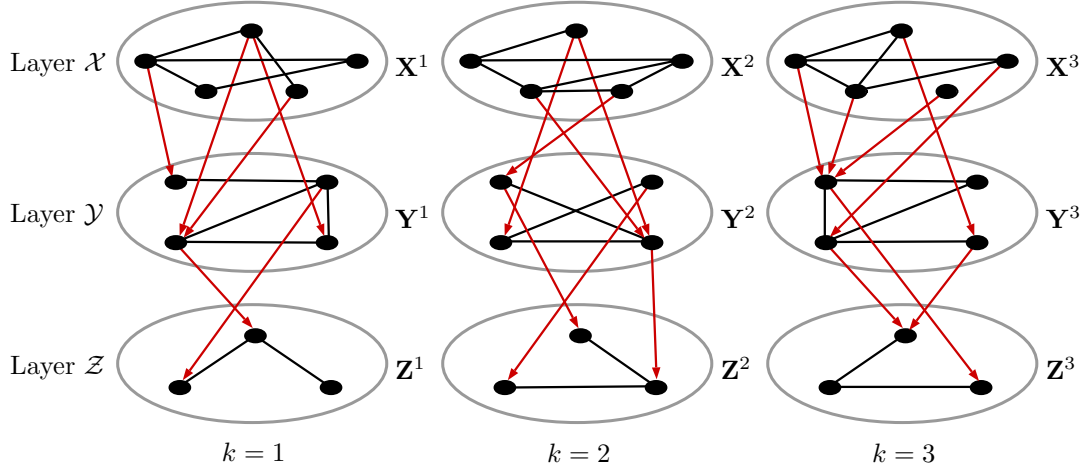


Figure 1: Multiple multilayer graphical models. The matrices  $(\mathbf{X}^k, \mathbf{Y}^k, \mathbf{Z}^k)$ ,  $k = 1, 2, 3$  indicate data for each layer and category  $k$ . Within-layer connections (black lines) are undirected, while between-layer connections (red lines) go from an upper layer to the successive lower layer. For each type of edges (i.e. within  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  and  $\mathcal{X} \rightarrow \mathcal{Y}, \mathcal{Y} \rightarrow \mathcal{Z}$ ), there are common edges across some or all  $k$ .

All the above approaches model either the horizontal or the vertical complexity in the full hierarchical structure of Figure 1. This means multiple related groups of heterogeneous datasets has to be modeled by analyzing all data in individual layers (i.e. models for  $\{\mathbf{X}^k\}, \{\mathbf{Y}^k\}, \{\mathbf{Z}^k\}$ ), and then separately analyzing individual hierarchies of datasets (i.e. separate models for  $(\mathbf{X}^k, \mathbf{Y}^k, \mathbf{Z}^k)$ ,  $k = 1, 2, 3$ ). Although a recent paper (Zhang et al., 2017) provides a model for the full structure in Figure 1 using penalized log-likelihoods, they do not give theoretical guarantees for the estimates, and limit the numerical examples to small feature dimensions ( $< 70$ ) only. Thus, a truly rigorous and scalable model for data integration is yet to be proposed.

While there has been some progress for parameter estimation in multilayer models, little is known about the asymptotic properties of resulting estimates. Current research on asymptotic distributions and testing procedures for estimates from high-dimensional problems has been limited to single-response regression using lasso (Javanmard and Montanari, 2014, 2018; van de Geer et al., 2014; Zhang and Zhang, 2014) or group lasso (Mitra and Zhang, 2016) penalties, and partial correlations of single (Cai and Liu, 2016) or multiple (Belilovsky et al., 2016; Liu, 2017) GGMs. From a systemic perspective, testing and identifying downstream interactions that differ across experimental conditions or disease subtypes can offer important insight on the underlying biological process (Li et al., 2015; Mao et al., 2017). In our framework, this can be done by developing a hypothesis testing procedure for entries in the within-layer regression matrices.

The contributions of this paper are two-fold. Firstly, we propose an integrative framework to conduct simultaneous inference for all parameters in multiple multilayer graphical models, essentially formalizing the structure in Figure 1. We decompose the multi-layer

problem into a series of two-layer problems, incorporate prior information on structural dependencies through imposing group structures on the model parameters, propose an estimation algorithm for them and derive theoretical properties of the estimators. Secondly, we obtain debiased versions of within-layer regression coefficients in this two-layer model, and derive their asymptotic distributions using estimates of model parameters that satisfy generic convergence guarantees. Consequently, we formulate a global test, as well as a simultaneous testing procedure that controls for False Discovery Rate (FDR) to detect important pairwise differences among directed edges between layers.

Our proposed framework for knowledge discovery from heterogeneous data sources is highly flexible. The group sparsity assumptions in our estimation technique can be replaced by other structural restrictions, for example low-rank or low-rank-plus-sparse, as and when deemed appropriate by the prior dependency assumptions across parameters. As long as the resulting estimates converge to corresponding true parameters at certain rates, they can be plugged into the testing methodology.

**Organization of paper.** We start with the model formulation in Section 2, then introduce our computational algorithm for a two-layer model, and derive theoretical convergence properties of the algorithm and resulting estimates. In section 3, we start by introducing the debiased versions of rows of the regression coefficient matrix estimates in our model, then use already computed parameter estimates that satisfy some general consistency conditions to obtain its asymptotic distribution. We then move on to pairwise testing, and use sparse estimates from our algorithm to propose a global test to detect overall differences in rows of the coefficient matrices, as well as a multiple testing procedure to detect elementwise differences and perform within-row thresholding of estimates in presence of moderate misspecification of the group sparsity structure. Section 4 is devoted to implementation of our methodology. We evaluate the performance of our estimation and testing procedure through several simulation settings, and give strategies to speed up the computational algorithm for high data dimensions. We conclude the paper with a discussion in Section 5. Proofs of all theoretical results, as well as some auxiliary results, are given in the appendix.

**Notations.** We denote scalars by small letters, vectors by bold small letters and matrices by bold capital letters. For any matrix  $\mathbf{A}$ ,  $(\mathbf{A})_{ij}$  denote its element in the  $(i, j)^{\text{th}}$  position. For  $a, b \in \mathbb{N}$ , we denote the set of all  $a \times b$  real matrices by  $\mathbb{M}(a, b)$ . For any positive integer  $c$ , define  $\mathcal{I}_c = \{1, \dots, c\}$ . For vectors  $\mathbf{v}$  and matrices  $\mathbf{M}$ ,  $\|\mathbf{v}\|$ ,  $\|\mathbf{v}\|_1$  or  $\|\mathbf{M}\|_1$  and  $\|\mathbf{v}\|_\infty$  or  $\|\mathbf{M}\|_\infty$  denote euclidean,  $\ell_1$  and  $\ell_\infty$  norms, respectively. The notation  $\text{supp}(\mathbf{A})$  indicates the non-zero edge set in a matrix (or vector)  $\mathbf{A}$ , i.e.  $\text{supp}(\mathbf{A}) = \{(i, j) : (\mathbf{A})_{ij} \neq 0\}$ . For positive real numbers  $A, B$  we write  $A \lesssim B$  if there exists  $c > 0$  independent of model parameters such that  $A \geq cB$ . We use the ‘:=’ notation to define a quantity for the first time.

## 2 The Joint Multiple Multilevel Estimation Framework

### 2.1 Formulation

Suppose there are  $K$  independent datasets, each pertaining to an  $M$ -layered Gaussian Graphical Model (GGM). The  $k^{\text{th}}$  model has the following structure:

$$\begin{aligned} \text{Layer 1-} & \quad \mathbb{D}_1^k = (D_{11}^k, \dots, D_{1p_1}^k) \sim \mathcal{N}(0, \Sigma_1^k); \quad k \in \mathcal{I}_K, \\ \text{Layer } m \text{ (} 1 < m \leq M \text{)-} & \quad \mathbb{D}_m^k = \mathbb{D}_{m-1}^k \mathbf{B}_m^k + \mathbb{E}_m^k, \text{ with } \mathbf{B}_m^k \in \mathbb{M}(p_{m-1}, p_m) \\ & \quad \text{and } \mathbb{E}_m^k = (E_{m1}^k, \dots, E_{mp_m}^k) \sim \mathcal{N}(0, \Sigma_m^k); \quad k \in \mathcal{I}_K. \end{aligned}$$

We assume known structured sparsity patterns, denoted by  $\mathcal{G}_m$  and  $\mathcal{H}_m$ , for the parameters of interest in the above model, i.e. the precision matrices  $\Omega_m^k := (\Sigma_m^k)^{-1}$  and the regression coefficient matrices  $\mathbf{B}_m^k$ , respectively. These patterns provide information on horizontal dependencies across  $k$  for the corresponding parameters, and our goal here is to leverage them to estimate the full hierarchical structure of the network- specifically to obtain the undirected edges inside nodes of a single layer, and the directed edges between two successive layers through jointly estimating  $\{\Omega_m^k\}$  and  $\{\mathbf{B}_m^k\}$ .

Consider now a two-layer model, which is a special case of the above model with  $M = 2$ :

$$\mathbb{X}^k = (X_1^k, \dots, X_p^k)^T \sim \mathcal{N}(0, \Sigma_x^k); \quad (2.1)$$

$$\mathbb{Y}^k = \mathbb{X}^k \mathbf{B}^k + \mathbb{E}^k; \quad \mathbb{E}^k = (E_1^k, \dots, E_p^k)^T \sim \mathcal{N}(0, \Sigma_y^k); \quad (2.2)$$

$$\mathbf{B}^k \in \mathbb{M}(p, q), \quad \Omega_x^k = (\Sigma_x^k)^{-1}; \quad \Omega_y^k = (\Sigma_y^k)^{-1}; \quad (2.3)$$

where we want to estimate  $\{(\Omega_x^k, \Omega_y^k, \mathbf{B}^k); k \in \mathcal{I}_K\}$  from data  $\mathcal{Z}^k = \{(\mathbf{Y}^k, \mathbf{X}^k); \mathbf{Y}^k \in \mathbb{M}(n, q), \mathbf{X}^k \in \mathbb{M}(n, p), k \in \mathcal{I}_K\}$  in presence of known grouping structures  $\mathcal{G}_x, \mathcal{G}_y, \mathcal{H}$  respectively and assuming  $n_k = n$  for all  $k \in \mathcal{I}_K$  for simplicity. We focus most of the theoretical discussion in the rest of the paper on jointly estimating  $\Omega_y := \{\Omega_y^k\}$  and  $\mathcal{B} := \{\mathbf{B}^k\}$ . This is because for  $M > 2$ , within-layer undirected edges of any  $m^{\text{th}}$  layer ( $m > 1$ ) and between-layer directed edges from the  $(m-1)^{\text{th}}$  layer to the  $m^{\text{th}}$  layer can be estimated from the corresponding data matrices in a similar fashion. On the other hand, parameters in the first layer are analogous to  $\Omega_x := \{\Omega_x^k\}$  that are dependent only on  $\{\mathbf{X}^k\}$ , so any method for joint estimation of multiple graphical models can be used to estimate them (e.g. Guo et al. (2011); Ma and Michailidis (2016)). This provides all building blocks for estimating the full hierarchical structure of our  $M$ -layered multiple GGMs.

### 2.2 Algorithm

We assume an element-wise group sparsity pattern over  $k$  for the precision matrices  $\Omega_x^k$ :

$$\mathcal{G}_x = \{\mathcal{G}_x^{ii'} : i \neq i'; i, i' \in \mathcal{I}_p\},$$

where each  $\mathcal{G}_x^{ii'}$  is a partition of  $\mathcal{I}_K$ . Subsequently we use the Joint Structural Estimation Method (JSEM) (Ma and Michailidis, 2016) to estimate  $\Omega_x$ , which first uses the group

structure given by  $\mathcal{G}_x$  in penalized nodewise regressions (Meinshausen and Bühlmann, 2006) to obtain neighborhood coefficients of each variable  $X_i, i \in \mathcal{I}_p$ , then fits a graphical lasso model over the combined support sets to obtain sparse estimates of the precision matrices:

$$\begin{aligned}\hat{\zeta}_i &= \arg \min_{\zeta_i} \left\{ \frac{1}{n} \sum_{k=1}^K \|\mathbf{X}_i^k - \mathbf{X}_{-i}^k \zeta_i^k\|^2 + \sum_{i' \leq i} \sum_{g \in \mathcal{G}_{ii'}^{x'}} \eta_n \|\zeta_{ii'}^{[g]}\| \right\} \\ \hat{E}_x^k &= \{(i, i') : 1 \leq i < i' \leq p, \hat{\zeta}_{ii'}^k \neq 0 \text{ OR } \hat{\zeta}_{i'i}^k \neq 0\} \\ \hat{\Omega}_x^k &= \arg \min_{\Omega_x^k \in \mathbb{S}_+(\hat{E}_x^k)} \left\{ \text{Tr}(\hat{\mathbf{S}}_x^k \Omega_x^k) - \log \det(\Omega_x^k) \right\}\end{aligned}\quad (2.4)$$

where  $\hat{\mathbf{S}}_x^k := (\mathbf{X}^k)^T \mathbf{X}^k / n_k$ .

For the precision matrices  $\Omega_y^k$  we assume an element-wise sparsity pattern  $\mathcal{G}_y$  defined in a similar manner as  $\mathcal{G}_x$ , while the sparsity pattern  $\mathcal{H}$  for  $\mathcal{B}$  is more general, each group  $h \in \mathcal{H}$  being defined as:

$$h = \{(\mathcal{S}_p, \mathcal{S}_q, \mathcal{S}_K) : \mathcal{S}_p \subseteq \mathcal{I}_p, \mathcal{S}_q \subseteq \mathcal{I}_q, \mathcal{S}_K \subseteq \mathcal{I}_K\}; \quad \bigcup_{h \in \mathcal{H}} h = \mathcal{I}_p \times \mathcal{I}_q \times \mathcal{I}_K$$

We obtain sparse estimates of  $\Omega_y$  and  $\mathcal{B}$  by solving the following group-penalized least square minimization problem:

$$\begin{aligned}\{\hat{\mathcal{B}}, \hat{\Theta}\} &= \arg \min_{\mathcal{B}, \Theta} \left\{ \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k - (\mathbf{Y}_{-j}^k - \mathbf{X}^k \mathbf{B}_{-j}^k) \boldsymbol{\theta}_j^k - \mathbf{X}^k \mathbf{B}_j^k\|^2 \right. \\ &\quad \left. + \sum_{j \neq j'} \sum_{g \in \mathcal{G}_{jj'}^y} \gamma_n \|\boldsymbol{\theta}_{jj'}^{[g]}\| + \sum_{h \in \mathcal{H}} \lambda_n \|\mathbf{B}^{[h]}\| \right\}\end{aligned}\quad (2.5)$$

$$\begin{aligned}\hat{E}_y^k &= \{(j, j') : 1 \leq j < j' \leq q, \hat{\boldsymbol{\theta}}_{jj'}^k \neq 0 \text{ OR } \hat{\boldsymbol{\theta}}_{j'j}^k \neq 0\} \\ \hat{\Omega}_y^k &= \arg \min_{\Omega_y^k \in \mathbb{S}_+(\hat{E}_y^k)} \left\{ \text{Tr}(\hat{\mathbf{S}}_y^k \Omega_y^k) - \log \det(\Omega_y^k) \right\}\end{aligned}\quad (2.6)$$

The outcome of a node in the lower layer is thus modeled using all other nodes in that layer *and* nodes in the immediate upper layer, with their effects quantified using  $\hat{\boldsymbol{\theta}}_j^k$  and  $\hat{\mathbf{B}}_j^k$ , respectively.

### 2.2.1 Alternating algorithm

The objective function in (2.5) is bi-convex, i.e. convex in  $\mathcal{B}$  for fixed  $\Theta$ , and vice-versa, but not jointly convex in  $\{\mathcal{B}, \Theta\}$ . Consequently, we use an alternating iterative algorithm to solve for  $\{\mathcal{B}, \Theta\}$  that minimizes (2.5) by iteratively cycling between  $\mathcal{B}$  and  $\Theta$ , i.e. holding one set of parameters fixed and solving for the other, then alternating until convergence.

Choice of initial values plays a crucial role in the performance of this alternating algorithm. We choose the initial values  $\{\widehat{\mathbf{B}}^{k(0)}\}$  by fitting separate lasso regression models for each column of the coefficient matrices:

$$\widehat{\mathbf{B}}_j^{k(0)} = \arg \min_{\mathbf{B}_j^k \in \mathbb{R}^p} \|\mathbf{Y}_j^k - \mathbf{X}^k \mathbf{B}_j^k\|^2 + \lambda_n \|\mathbf{B}_j^k\|_1; \quad j \in \mathcal{I}_q, k \in \mathcal{I}_K. \quad (2.7)$$

We obtain initial estimates of  $\Theta_j, j \in \mathcal{I}_q$  by performing group-penalized nodewise regression on the residuals  $\widehat{\mathbf{E}}^{k(0)} := \mathbf{Y}^k - \mathbf{X}^k \widehat{\mathbf{B}}^{k(0)}$ :

$$\widehat{\Theta}_j^{(0)} = \arg \min_{\Theta_j} \frac{1}{n} \sum_{k=1}^K \|\widehat{\mathbf{E}}_j^{k(0)} - \widehat{\mathbf{E}}_{-j}^{k(0)} \boldsymbol{\theta}_j^k\|^2 + \gamma_n \sum_{j \neq j'} \sum_{g \in \mathcal{G}_y^{jj'}} \|\boldsymbol{\theta}_{jj'}^{[g]}\|. \quad (2.8)$$

The steps of our full estimation procedure, which we call the Joint Multiple Multilayer Estimation (JMMLE) method, can thus be summarized in Algorithm 1.

**Algorithm 1.** (The JMMLE Algorithm)

1. Initialize  $\widehat{\mathcal{B}}$  using (2.7).
2. Initialize  $\widehat{\Theta}$  using (2.8).
3. Update  $\widehat{\mathcal{B}}$  as:

$$\widehat{\mathcal{B}}^{(t+1)} = \arg \min_{\substack{\mathbf{B}^k \in \mathbb{M}(p,q) \\ k \in \mathcal{I}_K}} \left\{ \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k - (\mathbf{Y}_{-j}^k - \mathbf{X}^k \mathbf{B}_{-j}^k) \widehat{\boldsymbol{\theta}}_j^{k(t)} - \mathbf{X}^k \mathbf{B}_j^k\|^2 + \lambda_n \sum_{h \in \mathcal{H}} \|\mathbf{B}^{[h]}\| \right\} \quad (2.9)$$

4. Obtain  $\widehat{\mathbf{E}}^{k(t+1)} := \mathbf{Y}^k - \mathbf{X}^k \mathbf{B}_j^{k(t)}, k \in \mathcal{I}_K$ . Update  $\widehat{\Theta}$  as:

$$\widehat{\Theta}_j^{(t+1)} = \arg \min_{\Theta_j \in \mathbb{M}(q-1,K)} \left\{ \frac{1}{n} \sum_{k=1}^K \|\widehat{\mathbf{E}}_j^{k(t+1)} - \widehat{\mathbf{E}}_{-j}^{k(t+1)} \boldsymbol{\theta}_j^k\|^2 + \gamma_n \sum_{j \neq j'} \sum_{g \in \mathcal{G}_y^{jj'}} \|\boldsymbol{\theta}_{jj'}^{[g]}\| \right\} \quad (2.10)$$

5. Continue till convergence.
6. Calculate  $\widehat{\Omega}_y^k, k \in \mathcal{I}_K$  using (2.6).

### 2.2.2 Tuning parameter selection

The nodewise regression step in the JSEM model (2.4) uses Bayesian Information Criterion (BIC) for tuning parameter selection. The step for updating  $\{\Theta\}$ , i.e. (2.10), in our JMMLE algorithm is analogous to this procedure, hence we use BIC to select the penalty parameter  $\gamma_n$ . In our setting the BIC for a given  $\gamma$  and fixed  $\mathcal{B}$  is given by:

$$\text{BIC}(\gamma; \mathcal{B}) = \text{Tr} \left( \mathbf{S}_y^k \widehat{\Omega}_{y,\gamma}^k \right) - \log \det \left( \widehat{\Omega}_{y,\gamma}^k \right) + \frac{\log n}{n} \sum_{k=1}^K |\widehat{E}_{y,\gamma}^k|$$

where  $\gamma$  in subscript indicates the corresponding quantity is calculated taking  $\gamma$  as the tuning parameter, and  $\mathbf{S}_y^k := (\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k)^T (\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k) / n$ . Every time  $\hat{\Theta}$  is updated in the JMMLE algorithm, we choose the optimal  $\gamma$  as the one with the smallest BIC over a fixed set of values  $\mathcal{C}_n$ . Thus for a fixed  $\lambda$ , our final choice of  $\gamma$  will be  $\gamma^*(\lambda) = \arg \min_{\gamma \in \mathcal{C}_n} \text{BIC}(\gamma; \hat{\mathcal{B}}_\lambda)$ .

We use the High-dimensional BIC (HBIC) to select the other tuning parameter,  $\lambda$ :

$$\begin{aligned} \text{HBIC}(\lambda; \Theta) = & \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k - (\mathbf{Y}_{-j}^k - \mathbf{X}^k \hat{\mathbf{B}}_{-j, \lambda}^k) \boldsymbol{\theta}_j^k - \mathbf{X}^k \hat{\mathbf{B}}_{j, \lambda}^k\|^2 + \\ & \log(\log n) \frac{\log(pq)}{n} \sum_{k=1}^K \left( \|\mathbf{B}^k\|_0 + |\hat{E}_{y, \gamma^*(\lambda)}^k| \right) \end{aligned}$$

We choose an optimal  $\lambda$  as the minimizer of HBIC by training multiple JMMLE models using Algorithm 1 over a finite set of values  $\lambda \in \mathcal{D}_n$ :  $\lambda^* = \arg \min_{\lambda \in \mathcal{D}_n} \text{HBIC}(\lambda, \hat{\Theta}_{\gamma^*(\lambda)})$ .

### 2.3 Properties of JMMLE estimators

We now provide theoretical results ensuring the convergence of our alternating algorithm, as well as the consistency of estimators obtained from the algorithm. We present statements of theorems in the main paper, giving detailed proofs and auxiliary results in the Appendix.

We introduce some notations that help establish the theorems that follow. Denote the true values of the parameters as  $\Omega_{x0} = \{\Omega_{x0}^k\}$ ,  $\Omega_{y0} = \{\Omega_{y0}^k\}$ ,  $\Theta_0 = \{\Theta_{0j}\}$ ,  $\mathcal{B}_0 = \{\mathcal{B}_0^k\}$ . Sparsity levels of individual true parameters are indicated by  $s_j := |\text{supp}(\Theta_{0j})|$ ,  $b_k := |\text{supp}(\mathcal{B}_0^k)|$ . Also define  $S := \sum_{j=1}^q s_j$ ,  $B := \sum_{k=1}^K b_k$ ,  $s := \max_{j \in \mathcal{I}_q} s_j$ .

Our first result establishes the convergence of Algorithm 1 for any fixed realization of  $\mathcal{X}$  and  $\mathcal{E}$ .

**Theorem 2.1.** *Suppose for any fixed  $(\mathcal{X}, \mathcal{E})$ , estimates in each iterate of Algorithm 1 are uniformly bounded by some quantity dependent on only  $p, q$  and  $n$ :*

$$\left\| (\hat{\mathcal{B}}^{(t)}, \hat{\Theta}_y^{(t)}) - (\mathcal{B}_0, \Theta_{y0}) \right\|_F \leq R(p, q, n); \quad t \geq 1 \quad (2.11)$$

*Then any limit point  $(\mathcal{B}^\infty, \Theta_y^\infty)$  of the algorithm is a stationary point of the objective function, i.e. a point where partial derivatives along all coordinates are non-negative.*

The next steps are to show that for random realizations of  $\mathcal{X}$  and  $\mathcal{E}$ ,

- (a) Successive iterates lie in this non-expanding ball around the true parameters,
- (b) The procedures in (2.7) and (2.8) ensure starting values that lie inside the same ball,

both with probability approaching 1 as  $(p, q, n) \rightarrow \infty$ .



To do so we break down the main problem into two subproblems. Take as  $\boldsymbol{\beta} = (\text{vec}(\mathbf{B}^1)^T, \dots, \text{vec}(\mathbf{B}^K)^T)^T$ : any subscript or superscript on  $\mathbf{B}$  being passed on to  $\boldsymbol{\beta}$ . Denote by  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\Theta}}$  the generic estimators given by

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{pqK}} \left\{ -2\boldsymbol{\beta}^T \hat{\boldsymbol{\gamma}} + \boldsymbol{\beta}^T \hat{\mathbf{\Gamma}} \boldsymbol{\beta} + \lambda \sum_{h \in \mathcal{H}} \|\boldsymbol{\beta}^{[h]}\| \right\} \quad (2.12)$$

$$\hat{\boldsymbol{\Theta}}_j = \arg \min_{\boldsymbol{\Theta}_j \in \mathbb{M}(q-1, K)} \left\{ \frac{1}{n} \sum_{k=1}^K \|\hat{\mathbf{E}}_j^k - \hat{\mathbf{E}}_{-j}^k \boldsymbol{\theta}_{jj'}^k\|^2 + \gamma \sum_{j \neq j'} \sum_{g \in \mathcal{G}_y^{jj'}} \|\boldsymbol{\theta}_{jj'}^{[g]}\| \right\}; \quad j \in \mathcal{I}_q \quad (2.13)$$

where

$$\hat{\mathbf{\Gamma}} = \begin{bmatrix} (\hat{\mathbf{T}}^1)^2 \otimes \frac{(\mathbf{X}^1)^T \mathbf{X}^1}{n} & & \\ & \ddots & \\ & & (\hat{\mathbf{T}}^K)^2 \otimes \frac{(\mathbf{X}^K)^T \mathbf{X}^K}{n} \end{bmatrix}; \quad \hat{\boldsymbol{\gamma}} = \begin{bmatrix} (\hat{\mathbf{T}}^1)^2 \otimes \frac{(\mathbf{X}^1)^T}{n} \\ \vdots \\ (\hat{\mathbf{T}}^K)^2 \otimes \frac{(\mathbf{X}^K)^T}{n} \end{bmatrix} \begin{bmatrix} \text{vec}(\mathbf{Y}^1) \\ \vdots \\ \text{vec}(\mathbf{Y}^K) \end{bmatrix}$$

with

$$\hat{T}_{jj'}^k = \begin{cases} 1 & \text{if } j = j' \\ -\hat{\theta}_{jj'}^k & \text{if } j \neq j' \end{cases} \quad (2.14)$$

It is easy to see that solving for  $\boldsymbol{\beta}$  in (2.5) given a fixed  $\hat{\boldsymbol{\Theta}}$  is equivalent to solving (2.12).

We assume the following conditions on the true parameter versions  $(\mathbf{T}_0^k)^2$ , defined from  $\boldsymbol{\Theta}_0$  similarly as (2.14):

**(E1)** The matrices  $(\mathbf{T}^k)^2, k \in \mathcal{I}_K$  are diagonally dominant, i.e.

$$|t_{0,jj}^k| > \sum_{j' \neq j} |t_{0,jj'}^k|$$

for  $j \in \mathcal{I}_q, k \in \mathcal{I}_K$ .

Now we are in a position to establish the estimation consistency for the solution of (2.12), given random  $(\mathcal{X}, \mathcal{E})$  and good enough estimators  $\hat{\boldsymbol{\Theta}}$ .

**Theorem 2.2.** Assume random  $(\mathcal{X}, \mathcal{E})$ , and fixed  $\hat{\boldsymbol{\Theta}}$  so that for  $j \in \mathcal{I}_q$ ,

$$\|\hat{\boldsymbol{\Theta}}_j - \boldsymbol{\Theta}_{0,j}\|_F \leq v_{\boldsymbol{\Theta}} = \eta_{\boldsymbol{\Theta}} \sqrt{\frac{\log q}{n}}$$

for some  $\eta_{\boldsymbol{\Theta}} > 0$  dependent on  $\boldsymbol{\Theta}$  only. Then, given the choice of tuning parameter

$$\lambda_n \geq 4\sqrt{|h_{\max}|} \mathbb{R}_0 \sqrt{\frac{\log(pq)}{n}}; \quad \mathbb{R}_0 := \max_{k \in \mathcal{I}_K} \mathbb{R}(v_{\boldsymbol{\Theta}}, \Sigma_x^k, \Sigma_y^k)$$

the following hold

$$\|\widehat{\beta} - \beta_0\|_1 \leq 48\sqrt{|h_{\max}|}B\lambda/\psi_* \quad (2.15)$$

$$\|\widehat{\beta} - \beta_0\| \leq 12\sqrt{B}\lambda/\psi_* \quad (2.16)$$

$$\sum_{h \in \mathcal{H}} \|\beta^{[h]} - \beta_0^{[h]}\| \leq 48B\lambda/\psi_* \quad (2.17)$$

$$(\widehat{\beta} - \beta_0)^T \widehat{\Gamma} (\widehat{\beta} - \beta_0) \leq 72B\lambda^2/\psi_* \quad (2.18)$$

with probability  $\geq 1 - 12c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n)$ , where  $|h_{\max}| = \max_{h \in \mathcal{H}} |h|$  and

$$\psi_* = \frac{1}{2} \min_k \left[ \Lambda_{\min}(\Sigma_{x_0}^k) \left( \min_j \psi_j^k - d_k v_{\Theta} \right) \right], \text{ with } \psi_j^k := t_{0,jj}^k - \sum_{j' \neq j} t_{0,jj'}^k$$

and  $d_k$  being the maximum degree of  $(\mathbf{T}_0^k)^2$ .

To prove an equivalent result for the solution of (2.13), as well as the consistency of the final estimates  $\widehat{\Omega}_y^k$  using their support sets, we need the following conditions.

**(E2)** For  $k \in \mathcal{I}_K$ ,  $\Omega_{y_0}^k$  is diagonally dominant, i.e.  $|\omega_{y_0,jj}| > \sum_{j' \neq j} |\omega_{y_0,jj'}|$  for  $j \in \mathcal{I}_q$ ;

**(E3)** There exist constants  $c_0, d_0$  such that for  $k \in \mathcal{I}_K$ ,

$$0 < 1/c_0 \leq \Lambda_{\min}(\Sigma_{y_0}^k) \leq \Lambda_{\max}(\Sigma_{y_0}^k) \leq 1/d_0 < \infty$$

Given these, we establish the required consistency results.

**Theorem 2.3.** Consider any deterministic  $\widehat{\mathbf{B}}$  that satisfy the following bound

$$\|\widehat{\mathbf{B}}^k - \mathbf{B}_0^k\|_1 \leq v_{\beta} = \eta_{\beta} \sqrt{\frac{\log(pq)}{n}}$$

Then, for sample size  $n \gtrsim \log(pq)$  the following hold:

(I) For the choice of tuning parameter  $\gamma \geq 4\sqrt{|g_{\max}|}\mathbb{Q}_0$ ,

$$\|\widehat{\Theta}_j - \Theta_{0,j}\|_F \leq 12\sqrt{s_j}\gamma/\psi \quad (2.19)$$

$$\sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\widehat{\theta}_{jj'}^{[g]} - \theta_{0,jj'}^{[g]}\| \leq 48s_j\gamma/\psi \quad (2.20)$$

(II) For the choice of tuning parameter  $\gamma = 4\sqrt{|g_{\max}|}\mathbb{Q}_0$ ,

$$\frac{1}{K} \sum_{k=1}^K \|\widehat{\Omega}_y^k - \Omega_y^k\|_F \leq O \left( \mathbb{Q}_0 \sqrt{\frac{|g_{\max}|S}{K}} \right) \quad (2.21)$$

both with probability  $\geq 1 - 1/p^{\tau_1-2} - 6c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n) - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$ , for some constants  $c_1, c_3, c_4 > 0, c_2, c_5 > 1$ .

Finally we ensure that the starting values are good enough.

**Theorem 2.4.** Consider the starting values as derived in (2.7) and (2.8). For sample size  $n \gtrsim \log(pq)$ , there exist constants  $d_1, d_2, d_3 > 0$  such that for

$$\lambda \geq 4d_2 \max_{k \in \mathcal{I}_K} \left\{ [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} \right\} \sqrt{\frac{\log(pq)}{n}}$$

we have  $\|\hat{\beta}^{(0)} - \beta_0\|_1 \leq 64B\lambda/\psi^*$  with probability  $\geq 1 - 6d_1 \exp(-(d_2^2 - 1) \log(pq)) - 2 \exp(d_3 n)$ . Further, for  $\gamma \geq 4\sqrt{|g_{\max}|} \mathbb{Q}_0$  we have

$$\begin{aligned} \|\hat{\Theta}_j^{(0)} - \Theta_{0,j}\|_F &\leq 12\sqrt{s_j}\gamma/\psi \\ \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\hat{\theta}_{jj'}^{[g](0)} - \theta_{0,jj'}^{[g]}\| &\leq 48s_j\gamma/\psi \end{aligned}$$

with probability  $\geq 1 - 1/p^{\tau_1-2} - 6c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n) - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$ .

Putting everything together, estimation consistency for the limit points of Algorithm 1 given our choice of starting values is immediate.

**Corollary 2.5.** Assume the conditions (E1)-(E3), and starting values  $\{\mathcal{B}^{(0)}, \Theta^{(0)}\}$  obtained using (2.7) and (2.8), respectively. Then, for random realizations of  $\mathcal{X}, \mathcal{E}$ ,

(I) For the choice of  $\lambda$

$$\lambda \geq 4 \max \left[ d_2 \max_{k \in \mathcal{I}_K} \left\{ [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} \right\}, \sqrt{|h_{\max}|} \mathbb{R}_0 \right] \sqrt{\frac{\log(pq)}{n}}$$

we have

$$\|\hat{\beta} - \beta_0\|_1 \leq \max \left\{ 48\sqrt{|h_{\max}|}, 64 \right\} \frac{B\lambda}{\psi_*}$$

with probability  $\geq \mathbf{tbd}$ .

(II) For  $\gamma \geq 4\sqrt{|g_{\max}|} \mathbb{Q}_0$ , (2.19) and (2.20) hold with probability  $\geq \mathbf{tbd}$ .

(III) For  $\gamma = 4\sqrt{|g_{\max}|} \mathbb{Q}_0$ , (2.21) holds with probability  $\geq \mathbf{tbd}$ .

### 3 Hypothesis testing in multilayer models

In this section, we lay out a framework for hypothesis testing in our proposed joint multilayer structure. Present literature in high-dimensional hypothesis testing either focuses on testing for similarities in the within-layer connections of single-layer networks (Cai and Liu, 2016; Liu, 2017), or coefficients of single response penalized regression (Mitra and Zhang, 2016; van de Geer et al., 2014; Zhang and Zhang, 2014). However, to our knowledge no method is available in the literature to perform testing for *between-layer* connections in a two-layer (or multilayer) setup.

There are two main challenges in doing the above: firstly the need to mitigate estimation the bias of estimators that are obtained from lasso and group lasso-based procedures and assumptions on the design matrix required for the same, and secondly the dependency among response nodes translating into the need for controlling False Discovery Rate (FDR) while simultaneously testing for such hypotheses. In Section 3.1 we propose a debiased estimator for rows of the coefficient matrix estimates  $\mathbf{B}^k$  that makes use of already computed (using JSEM) nodewise regression coefficients in the upper layer, and establish asymptotic properties of scaled version of them. Section 3.2 is devoted to pairwise testing, where we assume  $K = 2$ , and propose asymptotic global tests for detecting differential effects of a variable in the upper layer, as well as pairwise simultaneous tests for detecting elementwise difference in the coefficient matrices across  $k$ .

### 3.1 Debiased estimators and asymptotic normality

Zhang and Zhang (2014) proposed a debiasing procedure for lasso estimates and subsequently calculate confidence intervals for individual coefficients  $\beta_j$  in high-dimensional linear regression:  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{X} \in \mathbb{M}(n, p)$  and  $\epsilon_r \sim N(0, \sigma^2)$ ,  $r \in \mathcal{I}_n$  for some  $\sigma > 0$ . Given an initial lasso estimate  $\hat{\boldsymbol{\beta}}^{(init)} \in \mathbb{R}^p$  their debiased estimator was defined as:

$$\hat{\beta}_j^{(deb)} = \hat{\beta}_j^{(init)} + \frac{\mathbf{z}_j^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{(init)})}{\mathbf{z}_j^T \mathbf{x}_j}$$

where  $\mathbf{z}_j$  is the vector of residuals from  $\ell_1$ -penalized regression of  $\mathbf{x}_j$  on  $\mathbf{X}_{-j}$ . With centering around the true parameter value, say  $\beta_j^0$ , and proper scaling this has an asymptotic normal distribution:

$$\frac{\hat{\beta}_j^{(deb)} - \beta_j^0}{\|\mathbf{z}_j\|/|\mathbf{z}_j^T \mathbf{x}_j|} \sim N(0, \sigma^2)$$

Essentially, they obtain the debiasing factor for the  $j^{\text{th}}$  coefficient by taking residuals from the regularized regression and scale them using the projection of  $\mathbf{x}_j$  onto a space approximately orthogonal to it. Mitra and Zhang (2016) later generalized this idea to group lasso estimates. Further, van de Geer et al. (2014) and Javanmard and Montanari (2014) performed debiasing on the entire coefficient vectors.

We start off by defining debiased estimates for individual rows of the coefficient matrices  $\mathbf{B}^k$  in our two-layer model:

$$\hat{\mathbf{c}}_i^k = \hat{\mathbf{b}}_i^k + \frac{1}{nt_i^k} \left( \mathbf{X}_i^k - \mathbf{X}_{-i}^k \hat{\boldsymbol{\zeta}}_i^k \right)^T (\mathbf{Y}^k - \mathbf{X}^k \hat{\mathbf{B}}^k); \quad i \in \mathcal{I}_p, k \in \mathcal{I}_K \quad (3.1)$$

where  $t_i^k = (\mathbf{X}_i^k - \mathbf{X}_{-i}^k \hat{\boldsymbol{\zeta}}_i^k)^T \mathbf{X}_{-i}^k / n$ , and  $\hat{\boldsymbol{\zeta}}_i^k, \hat{\mathbf{B}}^k$  are *generic estimators* of the neighborhood coefficient matrices in the upper layer and within-layer coefficient matrices, respectively. By structure this is similar to the proposal of Zhang and Zhang (2014). However, as we see shortly, minimal conditions need to be imposed on the parameter estimates used in (3.1) for the asymptotic results based on a scaled version of the debiased estimator to go thorough, and they continue to hold for arbitrary sparsity patterns over  $k$  in all the parameters.

Present methods of debiasing coefficients from regularized regression require specific assumptions on the regularization structure of the main regression, as well as on how to calculate the debiasing factor. While [Zhang and Zhang \(2014\)](#), [Javanmard and Montanari \(2014\)](#) and [van de Geer et al. \(2014\)](#) work on coefficients from lasso regressions, [Mitra and Zhang \(2016\)](#) debias the coefficients of pre-specified groups in the coefficient vector from a group lasso. The current proposals for the debiasing factor available in the literature include nodewise lasso ([Zhang and Zhang, 2014](#)) and a variance minimization scheme with  $\ell_\infty$ -constraints ([Javanmard and Montanari, 2014](#)). In comparison, we only assume the following generic constraints on the parameter estimates used in our procedure.

**(T1)** For the upper layer neighborhood coefficients, the following holds for all  $k \in \mathcal{I}_K$ :

$$\|\hat{\boldsymbol{\zeta}}^k - \boldsymbol{\zeta}_0^k\|_1 \leq C_\zeta \sqrt{\frac{\log p}{n}}$$

where  $C_\zeta = O(1)$  depends only on the true values, i.e.  $\{\boldsymbol{\zeta}_0^k\}$ .

**(T2)** The lower layer precision matrix estimators satisfy for all  $k \in \mathcal{I}_K$

$$\left\| (\hat{\Omega}_y^k)^{1/2} - (\Omega_{y0}^k)^{-1/2} \right\|_\infty \leq C_\Omega \sqrt{\frac{\log q}{n}}$$

where  $C_\Omega = O(1)$  depends only on  $\Omega_{y0}$ .

**(T3)** For the regression coefficient matrices, the following holds for all  $k \in \mathcal{I}_K$ :

$$\|\hat{\mathbf{B}}^k - \mathbf{B}_0^k\|_1 \leq C_\beta \sqrt{\frac{\log(pq)}{n}}$$

where  $C_\beta = O(1)$  depends on  $\mathcal{B}$  only.

Given these conditions, the following theorem provides the asymptotic joint distribution of a scaled version of the debiased coefficients. A similar result for fixed design in the context of single-response linear regression can be found in the preprint by [Stucky and van de Geer \(2017\)](#). However they use nuclear norm as the loss function while obtaining the debiasing factors and use the resulting Karush-Kuhn-Tucker (KKT) conditions to derive their results, whereas we leverage bounds on generic parameter estimates combined with the sub-gaussianity of our design matrices.

**Theorem 3.1.** Define  $\hat{s}_i^k = \sqrt{\|\mathbf{X}_i^k - \mathbf{X}_{-i}^k \hat{\boldsymbol{\zeta}}_i^k\|^2/n}$ , and  $m_i^k = \sqrt{n t_i^k / \hat{s}_i^k}$ . Consider parameter estimates that satisfy conditions (T1)-(T3). Define the following:

$$\hat{\Omega}_y = \text{diag}(\hat{\Omega}_y^1, \dots, \hat{\Omega}_y^K)$$

$$\mathbf{M}_i = \text{diag}(m_i^1, \dots, m_i^K)$$

$$\mathbf{C}_i = \text{vec}(\hat{\mathbf{c}}_i^1, \dots, \hat{\mathbf{c}}_i^K)^T$$

$$\mathbf{D}_i = \text{vec}(\mathbf{b}_{0,i}^1, \dots, \mathbf{b}_{0,i}^K)^T$$

Then for sample size satisfying  $\log p = o(n^{1/2}), \log q = o(n^{1/2})$  we have

$$\widehat{\Omega}_y^{1/2} \mathbf{M}_i(\mathbf{C}_i - \mathbf{D}_i) \sim \mathcal{N}_{Kq}(\mathbf{0}, \mathbf{I}) + \mathbf{R}_n \quad (3.2)$$

where  $\|\mathbf{R}_n\|_\infty = o_P(1)$ .

### 3.2 Test formulation

Now we simply plug-in estimators from the JMMLE algorithm in Theorem 3.1. This is fairly straightforward. Condition (T1) is ensured by the JSEM penalized neighborhood estimators (immediate from Proposition A.1 in Ma and Michailidis (2016)), and a bound on total sparsity of the true coefficient matrices:  $B = o(n/\log(pq))$ , in conjunction with Corollary 2.5, ensures condition (T3), both with probability approaching 1 as  $(n, p, q) \rightarrow \infty$ . Finally, condition (T2) follows from Corollary 2.5 and Lemma A.9.

An asymptotic joint distribution of debiased versions of the JMMLE regression estimates is now immediate.

**Corollary 3.2.** *Consider the estimates  $\widehat{\mathbf{B}}$  and  $\widehat{\Omega}_y$  obtained from Algorithm 1, and upper layer neighborhood coefficients from solving the nodewise regression in (2.4). Suppose that  $\log(pq)/\sqrt{n} \rightarrow 0$ , and the sparsity condition  $B = o(n/\log(pq))$  is satisfied. Then, with the same notations in Theorem 3.1 we have*

$$\widehat{\Omega}_y^{1/2} \mathbf{M}_i(\mathbf{C}_i - \mathbf{D}_i) \sim \mathcal{N}_{Kq}(\mathbf{0}, \mathbf{I}) + \mathbf{R}_{1n} \quad (3.3)$$

where  $\|\mathbf{R}_{1n}\|_\infty = o_P(1)$ .

We are now ready to formulate asymptotic global and simultaneous testing procedures based on Corollary 3.2. In this paper, we restrict our attention to testing for pairwise differences only. Specifically, we set  $K = 2$ , and are interested in testing whether there are overall and elementwise differences between the coefficient vectors  $\mathbf{b}_{0i}^1$  and  $\mathbf{b}_{0i}^2$ .

When  $\mathbf{b}_{0i}^1 = \mathbf{b}_{0i}^2$ , it is immediate from Corollary 3.2 that a scaled version of the vector of estimated differences  $\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2$  follows a  $q$ -variate multinormal distribution. Consequently we formulate a global test for detecting differential overall downstream effect of the  $i^{\text{th}}$  covariate in the upper layer.

**Algorithm 2.** (Global test for  $H_0^i : \mathbf{b}_{0i}^1 = \mathbf{b}_{0i}^2$  at level  $\alpha, 0 < \alpha < 1$ )

1. Obtain the debiased estimators  $\widehat{\mathbf{c}}_i^1, \widehat{\mathbf{c}}_i^2$  using (3.1).
2. Calculate the test statistic

$$D_i = (\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2)^T \left( \frac{\widehat{\Sigma}_y^1}{(m_i^1)^2} + \frac{\widehat{\Sigma}_y^2}{(m_i^2)^2} \right)^{-1} (\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2)$$

where  $\widehat{\Sigma}_y^k = (\widehat{\Omega}_y^k)^{-1}, k = 1, 2$ .

3. Reject  $H_0^i$  if  $D_i \geq \chi_{q, 1-\alpha}^2$ .

Besides controlling the type-I error at a specified level  $\alpha$ , the above testing procedure maintains rate optimal power.

**Theorem 3.3.** Consider the global test given in Algorithm 2, performed using parameter estimates satisfying conditions (T1)-(T3). Say  $\boldsymbol{\delta} := \mathbf{b}_{0i}^1 - \mathbf{b}_{0i}^2$ . Then the power of the test is given by

$$K_q \left( \chi_{q,1-\alpha}^2 + \boldsymbol{\delta}^T \left( \frac{\Sigma_{y0}^1}{(m_{0i}^1)^2} + \frac{\Sigma_{y0}^2}{(m_{0i}^2)^2} + o(1) \right)^{-1} \boldsymbol{\delta} \right)$$

where  $K_q$  is the cumulative distribution function of the  $\chi_q^2$  distribution, and  $m_{0i}^k$  is the population version of  $m_i^k$ ,  $k = 1, 2$ . Consequently, for  $\|\boldsymbol{\delta}\|_\infty \geq \text{tbd}$ ,  $P(H_0^i \text{ is rejected}) \rightarrow 1$  as  $(n, p, q) \rightarrow \infty$ .

**proof tbd**

### 3.3 Control of False Discovery Rate

Given the null hypothesis is rejected, we consider the multiple testing problem of simultaneously testing for all entrywise differences, i.e. testing

$$H_0^{ij} : b_{0ij}^1 = b_{0ij}^2 \quad \text{vs.} \quad H_1^{ij} : b_{0ij}^1 \neq b_{0ij}^2$$

for all  $j \in \mathcal{I}_q$ . Here we use the test statistic

$$d_{ij} = \frac{\hat{c}_{ij}^1 - \hat{c}_{ij}^2}{\sqrt{\hat{\sigma}_{jj}^1/(m_i^1)^2 + \hat{\sigma}_{jj}^2/(m_i^2)^2}} \quad (3.4)$$

with  $\hat{\sigma}_{jj}^k$  being the  $j^{\text{th}}$  diagonal element of  $\hat{\Sigma}_y^k$ ,  $k = 1, 2$ .

For the purpose of simultaneous testing, we consider tests with a common rejection threshold  $\tau$ , i.e. for  $j \in \mathcal{I}_q$ ,  $H_0^{ij}$  is rejected if  $|d_{ij}| > \tau$ . We denote  $\mathcal{H}_0^i = \{j : b_{0,ij}^1 = b_{0,ij}^2\}$  and define the false discovery proportion (FDP) and false discovery rate (FDR) for these tests as follows:

$$FDP(\tau) = \frac{\sum_{j \in \mathcal{H}_0^i} \mathbb{I}(|d_{ij}| \geq \tau)}{\max \left\{ \sum_{j \in \mathcal{I}_q} \mathbb{I}(|d_{ij}| \geq \tau), 1 \right\}} \quad FDR(\tau) = \mathbb{E}[FDP(\tau)]$$

For a pre-specified level  $\alpha$ , we choose a threshold that ensures both FDP and FDR  $\leq \alpha$  using the Benjamini-Hochberg (BH) procedure. The procedure for FDR control is now given by Algorithm 3.

**Algorithm 3.** (Simultaneous tests for  $H_0^{ij} : b_{0ij}^1 = b_{0ij}^2$  at level  $\alpha$ ,  $0 < \alpha < 1$ )

1. Calculate the pairwise test statistics  $d_{ij}$  using (3) for  $j \in \mathcal{I}_q$ .
2. Obtain the threshold

$$\hat{\tau} = \inf \left\{ \tau \in \mathbb{R} : 1 - \Phi(\tau) \leq \frac{\alpha}{2q} \max \left( \sum_{j \in \mathcal{I}_q} \mathbb{I}(|d_{ij}| \geq \tau), 1 \right) \right\}$$

3. For  $j \in \mathcal{I}_q$ , reject  $H_0^{ij}$  if  $|d_{ij}| \geq \hat{\tau}$ .

For this procedure to maintain FDR and FDP asymptotically at a pre-specified level  $\alpha \in (0, 1)$ , we need some dependence conditions on true correlation matrices in the lower layer. Following [Liu and Shao \(2014\)](#), we introduce there are two types of dependency we consider:

**(D1)** Define  $r_{0,jj'}^k = \sigma_{0,jj'}^k / \sqrt{\sigma_{0,jj}^k \sigma_{0,j'j'}^k}$  for  $j, j' \in \mathcal{I}_q, k = 1, 2$ . Suppose there exists  $0 < r < 1$  such that  $\max_{1 \leq j < j' \leq q} |r_{0,jj'}^k| \leq r$ , and for every  $j \in \mathcal{I}_q$ ,

$$\sum_{j'=1}^q \mathbb{I} \left\{ |r_{0,jj'}^k| \geq \frac{1}{(\log q)^{2+\theta}} \right\} \leq O(q^\rho)$$

for some  $\theta > 0$  and  $0 < \rho < (1-r)/(1+r)$ .

**(D1\*)** Suppose there exists  $0 < r < 1$  such that  $\max_{1 \leq j < j' \leq q} |r_{0,jj'}^k| \leq r$ , and for every  $j \in \mathcal{I}_q$ ,

$$\sum_{j'=1}^q \mathbb{I} \left\{ |r_{0,jj'}^k| > 0 \right\} \leq O(q^\rho)$$

for some  $0 < \rho < (1-r)/(1+r)$ .

Originally proposed by [Liu and Shao \(2014\)](#), the above dependency conditions are meant to control the amount of correlation among the test statistics. Condition (D1) allows each variable to be highly correlated with at most  $O(q^\rho)$  other variables and weakly correlated with others, while (D1\*) limits the number of variables to have *any* correlation with it to  $O(q^\rho)$ . Note that (D1\*) is a stronger condition, and can be seen as the limiting condition of (D1) as  $q \rightarrow \infty$ .

**Theorem 3.4.** Suppose  $\mu_j = b_{0,ij}^1 - b_{0,ij}^2, \sigma_j^2 = n\mathbb{E}(\hat{\sigma}_{jj}^1/(m_i^1)^2 + \hat{\sigma}_{jj}^2/(m_i^2)^2)$ . Assume the following holds as  $(n, q) \rightarrow \infty$

$$\left| \left\{ j \in \mathcal{I}_q : |\mu_j/\sigma_j| \geq 4\sqrt{\log q/n} \right\} \right| \rightarrow \infty \quad (3.5)$$

Now Consider the conditions (D1) and (D1\*). If (D1) is satisfied, then the following holds when  $\log q = O(n^\xi), 0 < \xi < 3/23$ :

$$\frac{FDP(\hat{\tau})}{(|\mathcal{H}_0^i|/q)\alpha} \xrightarrow{P} 1; \quad \lim_{n,q \rightarrow \infty} \frac{FDR(\hat{\tau})}{(|\mathcal{H}_0^i|/q)\alpha} = 1 \quad (3.6)$$

Further, if (D1\*) is satisfied, then (3.6) holds for  $\log q = o(n^{1/3})$ .

The condition (3.5) is essential for FDR control in a diverging parameter space ([Liu, 2017; Liu and Shao, 2014](#)).

**Remark 1.** Following [Liu and Shao \(2014\)](#), a version of Algorithm 3, where the null distribution is calibrated using bootstrap instead of normal approximation, gives asymptotic FDR control under (D1\*) and  $\log q = o(n^{1/2})$ . We believe it is possible to obtain (3.6) under the weaker condition (D1) for  $\log q = o(n^{1/2})$  by extending the framework of [Liu \(2017\)](#) that performs multiple testing in multiple (single layer) GGMs, with the added advantage of being generalizable to the case of  $K > 2$ . However, this requires a significant amount of theoretical analysis, and we leave it for future research.



**Remark 2.** Based on the FDR control procedure in Algorithm 3, we can perform *within-row thresholding* in the matrices  $\widehat{\mathbf{B}}^k$  to tackle group misspecification.

$$\begin{aligned}\hat{\tau}_i^k &:= \inf \left\{ \tau \in \mathbb{R} : 1 - \Phi(\tau) \leq \frac{\alpha}{2q} \max \left( \sum_{j \in \mathcal{I}_q} \mathbb{I}(|b_{ij}^k| \geq \tau), 1 \right) \right\} \\ \hat{b}_{ij}^{k, \text{thr}} &= \hat{b}_{ij}^k \mathbb{I}(|\hat{b}_{ij}^k| \geq \hat{\tau}_i^k)\end{aligned}\tag{3.7}$$

Even without group misspecification, this helps identify directed edges between layers that have high nonzero values. Similar post-estimation thresholding results have been proposed in the context of multitask regression (Majumdar and Chatterjee, 2018; Obozinski et al., 2011) and neighborhood selection (Ma and Michailidis, 2016). However, our procedure is the first one to provide explicit guarantees on the amount of false discoveries while doing so.

## 4 Numerical performance

In this section, we evaluate the performance of our proposed JMMLE algorithm and the hypothesis testing framework in a two-layer simulation setup (Sections 4.1 and 4.2), and also introduce some computational techniques that significantly accelerates computation for high data dimensions (Section 4.3).

### 4.1 Simulation 1: estimation

As a first step towards obtaining a two-layer structure with horizontal (across  $k$ ) complexity and inter-layer directed edges, we generate the precision matrices  $\{\Omega_{x0}^k\}$  and  $\{\Omega_{y0}^k\}$  using a dependency structure across  $k$  that was first used in the simulation study of Ma and Michailidis (2016). We set  $K = 5$ , and set different shared sparsity patterns across  $k$  inside the lower  $p/2 \times p/2$  block of the upper layer precision matrices, and outside the block. In our notation, this means the following elementwise group structure:

$$\mathcal{G}_{x,ii'} = \begin{cases} \{(1, 2), (3, 4), 5\} & \text{if } i \leq p/2 \text{ or } j \leq p/2 \\ \{(1, 3, 5), (2, 4)\} & \text{otherwise} \end{cases}$$

The schematic in Figure 2 illustrates this structure. We set an off-diagonal element inside each of these common blocks (i.e.  $A, B, C$  and  $\alpha, \beta$  in the figure) to be non-zero with probability  $\pi_x \in \{5/p, 30/p\}$ , then generate the values of all non-zero elements independently from the uniform distribution in the interval  $[-1, 0.5] \cup [0.5, 1]$ . The precision matrices  $\Omega_{x0}^k$  are generated by putting together the corresponding common blocks, their positive definiteness ensured by setting all diagonal elements to be  $1 + |\Lambda_{\min}(\Omega_{x0}^k)|$ . We get elements in the covariance matrix as

$$\sigma_{x0,ii'}^k = (\omega_{x0,ii'}^k)^{-1} / \sqrt{(\omega_{x0,ii}^k)^{-1} (\omega_{x0,i'i'}^k)^{-1}},$$

and generate rows of  $\mathbf{X}^k$  independently from  $\mathcal{N}(0, \Sigma_{x0}^k)$ . We obtain  $\Sigma_{y0}^k$  and then  $\mathbf{E}^k$  using the same setup but with the number of variables being  $q$  and setting off-diagonal

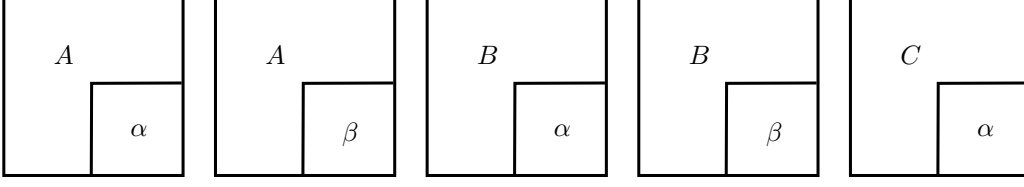


Figure 2: Shared sparsity patterns across  $k$  for the precision matrices  $\{\Omega_{x0}^k\}$  and  $\{\Omega_{y0}^k\}$

elements non-zero with probability  $\pi_y \in \{5/q, 30/q\}$ . To obtain the matrices  $\mathbf{B}_0^k$ , for a fixed  $(i, j)$ ,  $i \in \mathcal{I}_p, j \in \mathcal{I}_q$ , we set  $\mathbf{b}_{ij}^k$  non-zero across all  $k$  with probability  $\pi \in \{5/p, 5/q\}$ , generate the non-zero groups independently from  $\text{Unif}\{[-1, 0.5] \cup [0.5, 1]\}$ , and set  $\mathbf{Y}^k = \mathbf{X}^k \mathbf{B}_0^k + \mathbf{E}^k, k \in \mathcal{I}_K$ . Finally, we generate 50 such independent two-layer datasets for each of the following model settings:

- Set  $\pi_x = \pi = 5/p, \pi_y = 5/q$ , and  $(p, q, n) \in \{(60, 30, 100), (30, 60, 100), (200, 200, 150), (300, 300, 150)\}$ ;
- Set  $\pi_x = \pi = 30/p, \pi_y = 30/q$ , and  $(p, q, n) \in \{(200, 200, 100), (200, 200, 200)\}$ .

We use the following array of tuning parameters to train Algorithm 1:

$$\gamma \in \{0.3, 0.4, \dots, 1\} \sqrt{\frac{\log q}{n}}; \quad \lambda \in \{0.4, 0.6, \dots, 1.8\} \sqrt{\frac{\log p}{n}}$$

using the one-step version (Section 4.3) instead of the full algorithm to save computation time. We compare the performance of our joint estimation method with separate estimates of the parameters using the method of Lin et al. (2016), using the following performance metrics to evaluate estimates  $\tilde{\mathcal{B}} = \{\tilde{\mathbf{B}}^k\}$ :

- True positives-

$$\text{TP}(\tilde{\mathcal{B}}) = \frac{\sum_k |\text{supp}(\hat{\mathbf{B}}^k) \cup \text{supp}(\mathbf{B}_0^k)|}{\sum_k |\text{supp}(\mathbf{B}_0^k)|}$$

- True negatives-

$$\text{TN}(\tilde{\mathcal{B}}) = \frac{\sum_k |\text{supp}^c(\hat{\mathbf{B}}^k) \cup \text{supp}^c(\mathbf{B}_0^k)|}{\sum_k |\text{supp}^c(\mathbf{B}_0^k)|}$$

- Relative error in Frobenius norm-

$$\text{RF}(\tilde{\mathcal{B}}) = \sum_{k=1}^K \frac{\|\hat{\mathbf{B}}^k - \mathbf{B}_0^k\|_F}{\|\mathbf{B}_0^k\|_F}$$

We use the same metrics to evaluate the precision matrix estimates  $\tilde{\Omega}_y^k$  as well.

Tables 1 and 2 summarize the results. Joint estimation vastly outperforms separate estimation for  $\mathcal{B}_0$  across all metrics. JMMLE tends to be conservative for the estimation of  $\Omega_{y0}^k$ , producing sparse estimates and very high true negative proportions, although they are still far more accurate than those obtained from separate estimation, as evident by the low average RF scores.

$(\pi_x, \pi_y)$	$(p, q, n)$	Method	TP	TN	MCC	RF
$(5/p, 5/q)$	(60,30,100)	JMMLE	0.99 (0.002)	0.99 (0.01)		0.19 (0.019)
		Separate	0.95 (0.018)	0.99 (0.002)		0.27 (0.031)
	(30,60,100)	JMMLE	0.99 (0.004)	0.99 (0.006)		0.21 (0.014)
		Separate	0.66 (0.038)	0.99 (0.001)		0.59 (0.033)
	(200,200,150)	JMMLE	1.0 (0)	1.0 (0)		0.01 (0.005)
		Separate				
	(300,300,150)	JMMLE				
		Separate				
$(30/p, 30/q)$	(200,200,100)	JMMLE				
		Separate				
	(200,200,200)	JMMLE				
		Separate				

Table 1: Table of outputs for joint and separate estimation of regression matrices, giving empirical mean and standard deviation (in brackets) of each evaluation metric over 50 replications.

$(\pi_x, \pi_y)$	$(p, q, n)$	Method	TP	TN	MCC	RF
$(5/p, 5/q)$	(60,30,100)	JMMLE	0.66 (0.058)	0.95 (0.01)		0.33 (0.018)
		Separate	0.89 (0.018)	0.63 (0.014)		0.77 (0.044)
	(30,60,100)	JMMLE	0.58 (0.035)	0.98 (0.003)		0.32 (0.008)
		Separate	0.62 (0.027)	0.81 (0.007)		0.43 (0.011)
	(200,200,150)	JMMLE	0.39 (0.042)	0.99 (0.002)		0.30 (0.007)
		Separate				
	(300,300,150)	JMMLE				
		Separate				
$(30/p, 30/q)$	(200,200,100)	JMMLE				
		Separate				
	(200,200,200)	JMMLE				
		Separate				

Table 2: Table of outputs for joint and separate estimation of lower layer precision matrices over 50 replications.

$(\pi_x, \pi_y)$	$(p, q, n)$	TP( $\hat{\mathcal{B}}$ )	TN( $\hat{\mathcal{B}}$ )	MCC( $\hat{\mathcal{B}}$ )	RF( $\hat{\mathcal{B}}$ )
$(5/p, 5/q)$	(60,30,100)				
	(30,60,100)				
	(200,200,150)				
	(300,300,150)				
$(30/p, 30/q)$	(200,200,100)				
	(200,200,200)				
$(\pi_x, \pi_y)$	$(p, q, n)$	TP( $\hat{\Theta}$ )	TN( $\hat{\Theta}$ )	MCC( $\hat{\Theta}$ )	RF( $\hat{\Theta}$ )
$(5/p, 5/q)$	(60,30,100)				
	(30,60,100)				
	(200,200,150)				
	(300,300,150)				
$(30/p, 30/q)$	(200,200,100)				
	(200,200,200)				

Table 3: Table of outputs for joint estimation in presence of group misspecification.

#### 4.1.1 Effect of heterogeneity

We repeat the above setups to check the performance of JMMLE in presence of within-group misspecification. For this we set individual elements in a non-zero group  $\{b_{ij}^k, k \in \mathcal{I}_K\}$  to be non-zero with probability 0.2, then pass JMMLE estimates of  $\mathbf{B}_0^k$  through the FDR controlling thresholds as given in (3.7). The results are summarized in ??.

## 4.2 Simulation 2: testing

We slightly change our data generating model to evaluate our proposed global testing and FDR control procedure. We set  $K = 2$ , then generate the  $\mathbf{B}_0^1$  by first randomly assigning each of its element to be non-zero with probability  $\pi$ , then drawing values of those elements from  $\text{Unif}\{[-1, -0.5] \cup [0.5, 1]\}$  independently. After this we generate a matrix of differences  $\mathbf{D}$ , with  $(\mathbf{D})_{ij}, i \in \mathcal{I}_p, j \in \mathcal{I}_q$  taking values -1, 1, 0 with probabilities 0.1, 0.1 and 0.8, respectively. Finally we set  $\mathbf{B}_0^2 = \mathbf{B}_0^1 + \mathbf{D}$ . We set the same sparsity structures for the pairs of precision matrices  $\{\Omega_{x0}^1, \Omega_{x0}^2\}$  and  $\{\Omega_{y0}^1, \Omega_{y0}^2\}$ . We use 50 replications of the above setup to calculate empirical power of global tests, as well as empirical power and FDR of simultaneous tests, while to get size of global tests we use JMMLE estimators from a separate set of data generated setting all elements of  $\mathbf{D}$  to 0. The type-I error of global tests is taken as 0.05, while FDR is set at 0.2 while calculating the respective thresholds.

Table 4 reports the empirical mean and standard deviations (in brackets) of all relevant quantities. We report outputs for all combinations of data dimensions and sparsity used in Section 4.1, and also for increased sample sizes in each setting until a satisfactory FDR is reached. As expected, higher sample sizes result in increased power for both global and simultaneous tests, and decreased size and FDR for all but one ( $p = 30, q = 60$ ) of the settings.

$(\pi_x, \pi_y)$	$(p, q)$	$n$	Global test		Simultaneous tests	
			Power	Size	Power	FDR
$(5/p, 5/q)$	$(60, 30)$	100	0.977 (0.018)	0.058 (0.035)	0.937 (0.021)	0.237 (0.028)
		200	0.987 (0.016)	0.046 (0.032)	0.968 (0.013)	0.218 (0.032)
	$(30, 60)$	100	0.985 (0.018)	0.097 (0.069)	0.925 (0.022)	0.24 (0.034)
		200	0.990 (0.02)	0.119 (0.059)	0.958 (0.024)	0.245 (0.041)
	$(200, 200)$	150	0.987 (0.005)	0.004 (0.004)	0.841 (0.13)	0.213 (0.007)
	$(300, 300)$	150	0.988 (0.002)	0.002 (0.003)	0.546 (0.035)	0.347 (0.017)
		300	0.998 (0.003)	0.000 (0.001)	0.989 (0.003)	0.117 (0.006)
$(30/p, 30/q)$	$(200, 200)$	100	0.994 (0.005)	0.262 (0.06)	0.479 (0.01)	0.557 (0.006)
		200	0.998 (0.004)	0.020 (0.01)	0.962 (0.003)	0.266 (0.007)
		300	0.999 (0.002)	0.011 (0.008)	0.990 (0.004)	0.185 (0.009)

Table 4: Table of outputs for hypothesis testing.

### 4.3 Computation

We now discuss some observations and strategies that speed up the JMMLE algorithm and reduces computations time significantly, especially for higher number of features in either layer.

**Block update and refit  $\mathbf{B}^k$  in each iteration.** Similar to the case of  $K = 1$  (Lin et al., 2016), we use block coordinate descent *within* each  $\mathbf{B}^k$ . This means instead of the full update step (2.9) we perform the following steps in each iteration to speed up convergence:

$$\left\{ \hat{\mathbf{b}}_j^{k(t+1)} \right\}_{k=1}^K = \arg \min_{\substack{\mathbf{b}_j^k \in \mathbb{R}^p \\ k \in \mathcal{I}_K}} \left\{ \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k + \mathbf{r}_j^{k(t)} - \mathbf{X}^k \mathbf{b}_j^k\|^2 + \lambda \sum_{h \in \mathcal{H}} \|\mathbf{b}_j^{[h]}\| \right\}$$

where  $\mathbf{r}_1^{k(t)} = \hat{\mathbf{E}}_{-1}^{k(t)} \hat{\boldsymbol{\theta}}_1^{k(t)}$ , and

$$\mathbf{r}_j^{k(t)} = \sum_{j'=1}^{j-1} \hat{\mathbf{e}}_j^{k(t+1)} \hat{\theta}_{jj'}^{k(t)} + \sum_{j'=j+1}^q \hat{\mathbf{e}}_j^{k(t)} \hat{\theta}_{jj'}^{k(t)}$$

for  $j \geq 2$ . Further, when starting from the initializer of the coefficient matrix given in (2.7), the support set of coefficient estimates becomes constant after only a few ( $\sim 10$ ) iterations of our algorithm, after which it refines the values inside the same support until overall convergence. This process speeds up significantly if a refitting step is added *in each iteration* after the matrices  $\hat{\mathbf{B}}^k$  are updated:

$$\left\{ \tilde{\mathbf{b}}_j^{k(t+1)} \right\}_{k=1}^K = \arg \min_{\substack{\mathbf{b}_j^k \in \mathbb{R}^p \\ k \in \mathcal{I}_K}} \left\{ \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k + \mathbf{r}_j^{k(t)} - \mathbf{X}^k \mathbf{b}_j^k\|^2 + \lambda \sum_{h \in \mathcal{H}} \|\mathbf{B}_{-j}^{[h]}\| \right\};$$

$$\hat{\mathbf{b}}_j^{k(t+1)} = \left[ (\mathbf{X}_{S_{jk}}^k)^T (\mathbf{X}_{S_{jk}}^k) \right]^{-1} (\mathbf{X}_{S_{jk}}^k)^T \mathbf{Y}_j^k$$

$(p, q, n)$	Method	TP( $\hat{\mathcal{B}}$ )	TN( $\hat{\mathcal{B}}$ )	MCC( $\hat{\mathcal{B}}$ )	RF( $\hat{\mathcal{B}}$ )
(60,30,100)	Full	0.999 (0.002)	0.992 (0.009)		0.195 (0.021)
	One step	0.999 (0.002)	0.993 (0.01)		0.190 (0.019)
(30,60,100)	Full	0.997 (0.004)	0.986 (0.007)		0.205 (0.014)
	One step	0.996 (0.004)	0.988 (0.006)		0.206 (0.014)
$(p, q, n)$	Method	TP( $\hat{\Theta}$ )	TN( $\hat{\Theta}$ )	MCC( $\hat{\Theta}$ )	RF( $\hat{\Theta}$ )
(60,30,100)	Full	0.671 (0.052)	0.949 (0.01)		0.327 (0.015)
	One step	0.663 (0.058)	0.95 (0.01)		0.328 (0.018)
(30,60,100)	Full	0.58 (0.039)	0.982 (0.003)		0.32 (0.009)
	One step	0.577 (0.035)	0.981 (0.003)		0.321 (0.008)

Table 5: Comparison of evaluation metrics for full and one-step versions of the JMMLE algorithm.

where  $\mathcal{S}_{jk} = \text{supp}(\tilde{\mathbf{b}}_j^{k(t+1)})$ .

**One-step estimator.** Algorithm 1, even after the above modifications, is computation-intensive. The reason behind this is the full tuning and updating of the lower layer neighborhood estimates  $\{\hat{\Theta}_j\}$  in each iteration. In practice, the algorithm speeds up significantly without compromising on estimation accuracy if we dispense of the  $\Theta_j$  updation step in all but the last iteration. More precisely, we consider the following one-step version of the original algorithm.

**Algorithm 4.** (The one-step JMMLE Algorithm)

1. Initialize  $\hat{\mathcal{B}}$  using (2.7).
2. Initialize  $\hat{\Theta}$  using (2.8).
3. Update  $\hat{\mathcal{B}}$  as:

$$\hat{\mathcal{B}}^{(t+1)} = \arg \min_{\substack{\mathbf{B}^k \in \mathbb{M}(p,q) \\ k \in \mathcal{I}_K}} \left\{ \frac{1}{n} \sum_{j=1}^q \sum_{k=1}^K \|\mathbf{Y}_j^k - (\mathbf{Y}_{-j}^k - \mathbf{X}^k \mathbf{B}_{-j}^k) \hat{\boldsymbol{\theta}}_j^{k(0)} - \mathbf{X}^k \mathbf{B}_j^k\|^2 + \lambda_n \sum_{h \in \mathcal{H}} \|\mathbf{B}^{[h]}\| \right\}$$

4. Continue till convergence to obtain  $\hat{\mathcal{B}} = \{\hat{\mathbf{B}}^k\}$ .
5. Obtain  $\hat{\mathbf{E}}^k := \mathbf{Y}^k - \mathbf{X}^k \hat{\mathbf{B}}^k, k \in \mathcal{I}_K$ . Update  $\hat{\Theta}$  as:

$$\hat{\Theta}_j = \arg \min_{\Theta_j \in \mathbb{M}(q-1,K)} \left\{ \frac{1}{n} \sum_{k=1}^K \|\hat{\mathbf{E}}_j^k - \hat{\mathbf{E}}_{-j}^k \boldsymbol{\theta}_j^k\|^2 + \gamma \sum_{j \neq j'} \sum_{g \in \mathcal{G}_{jj'}} \|\boldsymbol{\theta}_{jj'}^{[g]}\| \right\}$$

6. Calculate  $\hat{\Omega}_y^k, k \in \mathcal{I}_K$  using (2.6).

Table 5 compares performances the full algorithm and the one-step version for the two data settings with smaller feature dimensions. The performances are indistinguishable across metrics, but the one-step algorithm saves computation time in orders of magnitude (Table 6).

$(p, q, n)$	Method	Comp. time (min)
(60,30,100)	Full	
	One-step	
(30,60,100)	Full	19.8
	One-step	3.1

Table 6: Comparison of computation times (averaged over 50 replications) for full and one-step versions of the JMMLE algorithm.

## 5 Discussion

Our work introduces an integrative framework for knowledge discovery in multiple multi-layer Gaussian Graphical Models. We exploit *a priori* known structural similarities across parameters of the multiple models to achieve estimation gains compared to separate estimation. More importantly, we derive results on the asymptotic distributions of generic estimates of the multiple regression coefficient matrices in this complex setup, and perform global and simultaneous testing for pairwise differences within the between-layer edges.

Our hypothesis testing framework has two immediate extensions.

- (I) In a recent work, [Liu \(2017\)](#) proposed a framework to test for structural similarities and differences across multiple *single layer* GGMs. For  $K$  GGMs with precision matrices  $\Omega^k = (\omega_{ii'})_{i,i' \in \mathcal{I}_p}$ , they test for the partial correlation coefficients  $\rho_{ii'}^{(k)} = -\omega_{ii'}^{(k)} / \sqrt{\omega_{ii}^{(k)} \omega_{i'i'}^{(k)}}$  using residuals from  $pK$  separate penalized neighborhood regressions, one for each variable of each GGM. To incorporate structured sparsity across  $k$ , our simultaneous regression techniques for all neighborhood coefficients (i.e. (2.4) and (2.13)) can be used instead to perform testing on the between-layer edges. Theoretical properties of this procedure can be derived using results in [Liu \(2017\)](#), possibly with adjustments for our neighborhood estimates to adhere to the rate conditions for the constants  $a_{n1}, a_{n2}$  therein to account for a diverging  $(p, q, n)$  setup.
- (II) For  $K \geq 2$ , detection of the following sets of inter-layer edges can be scientifically significant:

$$\begin{aligned}
\mathcal{B}_1 &= \left\{ (i, j) : \sum_{1 \leq k < k' \leq K} \left( b_{0,ij}^k - b_{0,ij}^{k'} \right)^2 > 0; i \in \mathcal{I}_p, j \in \mathcal{I}_q \right\} \\
\mathcal{B}_2 &= \{ (i, j) : b_{0,ij}^1 = \cdots, b_{0,ij}^K \neq 0 \} \\
\mathcal{B}_3 &= \{ (i, j) : b_{0,ij}^1 = \cdots, b_{0,ij}^K = 0 \}
\end{aligned}$$

e.g. detection of gene-protein interactions that are present, but may have different or same weights across  $k$  ( $\mathcal{B}_1$  and  $\mathcal{B}_2$ , respectively), and that are absent for all  $k$  ( $\mathcal{B}_3$ ). The asymptotic result in Theorem 3.1 continues to hold in this situation, and an extension of the global test (Algorithm 2) is immediate. However, extending the FDR control procedure requires a technically more detailed approach.

The strength of our proposed debiased estimators (3.1) is that only generic estimators of relevant model parameters that satisfy general rate conditions are necessary for it to have a valid asymptotic distribution. This translates to a high degree of flexibility while choosing the method of estimation. Our formulation based on sparsity assumptions (Section 2.2) is only one such way to obtain the necessary estimates. Sparsity may not be an assumption that is required or even valid in complex hierarchical structures from different domains of application. For different two-layer pairs in such multilayer setups, low-rank, group-sparse or sparse methods (or a combination thereof) can be plugged into our alternating algorithm. Results parallel to those in Section 2.3 need to be established for the corresponding estimators. However, as long as these estimators adhere to the convergence conditions (T1)-(T3), Theorem 3.1 can be used to derive the asymptotic distributions of between-layer edges.



## Appendix

### A Proof of main results

**need modification** To prove the results in this section, we use a reparametrization of the neighborhood coefficients at the lower level. Specifically, notice that for  $j \in \mathcal{I}_q, k \in \mathcal{I}_K$ , the corresponding summand in  $f(\mathcal{Y}, \mathcal{X}, \mathcal{B}, \Theta)$  can be rearranged as

$$\begin{aligned} \|\mathbf{Y}_j^k - \mathbf{X}^k \mathbf{B}_j^k - (\mathbf{Y}_{-j}^k - \mathbf{X}^k \mathbf{B}_{-j}^k) \boldsymbol{\theta}_j^k\|^2 &= \|\mathbf{Y}_j^k - \mathbf{Y}_{-j}^k \boldsymbol{\theta}_j^k - (\mathbf{X}^k \mathbf{B}_j^k - \mathbf{X}^k \mathbf{B}_{-j}^k) \boldsymbol{\theta}_j^k\|^2 \\ &= \|(\mathbf{Y} - \mathbf{X} \mathbf{B}) \mathbf{T}_j^k\|^2 \end{aligned}$$

where

$$T_{jj'}^k = \begin{cases} 1 & \text{if } j = j' \\ -\theta_{jj'}^k & \text{if } j \neq j' \end{cases}$$

Thus, with  $\mathbf{T}^k := (\mathbf{T}_j^k)_{j \in \mathcal{I}_q}$ , we have

$$f(\mathcal{Y}, \mathcal{X}, \mathcal{B}, \Theta) = \frac{1}{n} \sum_{j=1}^p \sum_{k=1}^K \|(\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k) \mathbf{T}_j^k\|^2 = \frac{1}{n} \sum_{k=1}^K \|\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k\|_{\mathbf{T}^k}^2 = \sum_{k=1}^K \text{Tr}(\mathbf{S}^k (\mathbf{T}^k)^2)$$

where  $\mathbf{S}^k = (1/n)(\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k)(\mathbf{Y}^k - \mathbf{X}^k \mathbf{B}^k)^T$  is the sample covariance matrix.

**Theorem A.1.** Assume fixed  $\mathcal{X}, \mathcal{E}$  and deterministic  $\widehat{\mathbf{B}} = \{\widehat{\mathbf{B}}^k\}$ . Also for  $k = 1, \dots, K$ ,

(T1)  $\|\widehat{\mathbf{B}}^k - \mathbf{B}_0^k\|_1 \leq v_\beta$ , where  $v_\beta = \eta_\beta \sqrt{\frac{\log(pq)}{n}}$  with  $\eta_\beta \geq 0$  depending on  $\mathcal{B}$  only;

(T2) Denote  $\widehat{\mathbf{E}}^k = \mathbf{Y}^k - \mathbf{X}^k \widehat{\mathbf{B}}^k, k \in \mathcal{I}_K$ . Then for all  $j \in \mathcal{I}_q$ ,

$$\frac{1}{n} \left\| (\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_\infty \leq \mathbb{Q}(v_\beta, \Sigma_x^k, \Sigma_y^k)$$

where  $\mathbb{Q}(v_\beta, \Sigma_x^k, \Sigma_y^k)$  is a  $O(\sqrt{\log(pq)/n})$  deterministic function of  $\mathcal{B}, \Sigma_x^k$  and  $\Sigma_y^k$ .

(T3) Denote  $\widehat{\mathbf{S}}^k = (\widehat{\mathbf{E}}^k)^T \widehat{\mathbf{E}}^k / n$ . Then  $\widehat{\mathbf{S}}^k \sim RE(\psi^k, \phi^k)$  with  $Kq\phi \leq \psi/2$  where  $\psi = \min_k \psi^k, \phi = \max_k \phi^k$ ;

(T4) Assumption (A2) holds for  $\Sigma_y^k$ .

Then, given the choice of tuning parameter

$$\gamma_n = 4\sqrt{|g_{\max}|} \mathbb{Q}_0; \quad \mathbb{Q}_0 := \max_{k \in \mathcal{I}_K} \mathbb{Q}(v_\beta, \Sigma_x^k, \Sigma_y^k)$$

the following holds

$$\frac{1}{K} \sum_{k=1}^K \|\widehat{\Omega}_y^k - \Omega_y^k\|_F \leq O\left(\mathbb{Q}_0 \sqrt{\frac{|g_{\max}|S}{K}}\right)$$

where  $|g_{\max}|$  is the maximum group size.

When  $\mathcal{X}$  and  $\mathcal{E}$  are random, the following propositions ensures that conditions (T2) and (T3) hold with probabilities approaching to 1.

**Proposition A.2.** Consider deterministic  $\widehat{\mathcal{B}}$  satisfying assumption (T1). Then for sample size  $n \gtrsim \log(pq)$  and  $k \in \mathcal{I}_K$ ,

1.  $\widehat{\mathbf{S}}^k$  satisfies the RE condition:  $\widehat{\mathbf{S}}^k \sim RE(\psi^k, \phi^k)$ , where

$$\psi^k = \frac{\Lambda_{\min}(\Sigma_x^k)}{2}; \quad \phi^k = \frac{\psi^k \log p}{n} + 2v_\beta c_2 [\Lambda_{\max}(\Sigma_x^k) \Lambda_{\max}(\Sigma_y^k)]^{1/2} \sqrt{\frac{\log(pq)}{n}}$$

with probability  $\geq 1 - 6c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n)$ ,  $c_1, c_3 > 0, c_2 > 1$ .

2. The following deviation bound is satisfied for any  $j \in \mathcal{I}_q$

$$\left\| \frac{1}{n} (\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_\infty \leq \mathbb{Q} \left( v_\beta, \Sigma_x^k, \Sigma_y^k \right)$$

with probability  $\geq 1 - 1/p^{\tau_1 - 2} - 6c_1 \exp[-(c_2^2 - 1) \log(pq)] - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$ ,  $c_4 > 0, c_5 > 1$ , where

$$\begin{aligned} \mathbb{Q} \left( v_\beta, \Sigma_x^k, \Sigma_y^k \right) &= 2v_\beta^2 V_x^k + 4v_\beta c_2 [\Lambda_{\max}(\Sigma_x^k) \Lambda_{\max}(\Sigma_y^k)]^{1/2} \sqrt{\frac{\log(pq)}{n}} + \\ &\quad c_5 \left[ \Lambda_{\max}(\Sigma_{y,-j}^k) \sigma_{y,j,-j}^k \right]^{1/2} \sqrt{\frac{\log q}{n}} \end{aligned}$$

with  $\sigma_{y,j,-j}^k = \mathbb{V}(E_j - \mathbb{E}_{-j} \boldsymbol{\theta}_{0,j})$ , and

$$V_x^k = \sqrt{\frac{\log 4 + \tau_1 \log p}{c_x^k n}} + \max_i \sigma_{x,ii}^k; \quad c_x^k = \left[ 128(1 + 4\Lambda_{\max}(\Sigma_x))^2 \max_i (\sigma_{x,ii})^2 \right]^{-1}$$

The error bounds for  $\widehat{\Omega}_y^k, k \in \mathcal{I}_K$  follow immediately from the above two results.

**Corollary A.3.** Consider any deterministic  $\widehat{\mathcal{B}}$  that satisfy the following bound

$$\|\widehat{\mathbf{B}}^k - \mathbf{B}_0^k\|_1 \leq v_\beta = \eta_\beta \sqrt{\frac{\log(pq)}{n}}$$

Then, for sample size  $n \gtrsim \log(pq)$  and choice of tuning parameter  $\gamma_n = 4\sqrt{|g_{\max}|} \mathbb{Q}_0$ , there exist constants  $c_1, c_3, c_4 > 0, c_2, c_5 > 1$  such that the following holds

$$\frac{1}{K} \sum_{k=1}^K \|\widehat{\Omega}_y^k - \Omega_y^k\|_F \leq O \left( \mathbb{Q}_0 \sqrt{\frac{|g_{\max}| S}{K}} \right) \quad (\text{A.1})$$

with probability  $\geq 1 - 1/p^{\tau_1 - 2} - 6c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n) - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$ .

### Discuss tighter bound compared to vanilla JSEM

After providing the error bounds for solutions to the subproblem (2.13), we concentrate on the subproblem (2.12). Following a similar strategy, we first get error bounds for  $\widehat{\beta}$  assuming everything else fixed.

**Theorem A.4.** Assume fixed  $\mathcal{X}, \mathcal{E}$ , and deterministic  $\hat{\Theta} = \{\hat{\Theta}_j\}$ , so that for  $j \in \mathcal{I}_q$ ,

(B1)  $\|\hat{\Theta}_j - \Theta_{0,j}\|_F \leq v_\Theta \sqrt{\frac{\log q}{n}}$  for some  $v_\Theta$  dependent on  $\Theta$ .

(B2) Denote  $\hat{\Gamma}^k = (\hat{\mathbf{T}}^k)^2 \otimes (\mathbf{X}^k)^T \mathbf{X}^k / n$ ,  $\hat{\gamma}^k = (\hat{\mathbf{T}}^k)^2 \otimes (\mathbf{X}^k)^T \mathbf{Y}^k / n$ . Then the deviation bound holds:

$$\left\| \hat{\gamma}^k - \hat{\Gamma}^k \beta_0 \right\|_\infty \leq \mathbb{R}(v_\Theta, \Sigma_x^k, \Sigma_y^k) \sqrt{\frac{\log(pq)}{n}}$$

where  $\mathbb{R}(v_\Theta, \Sigma_x^k, \Sigma_y^k)$  is a  $O(1)$  deterministic function of  $\Theta, \Sigma_x^k$  and  $\Sigma_y^k$ .

(B3)  $\hat{\Gamma} \sim RE(\psi_*, \phi_*)$  with  $Kpq\phi_* \leq \psi_*/2$ .

Then, given the choice of tuning parameter

$$\lambda_n \geq 4\sqrt{|h_{\max}|} \mathbb{R}_0 \sqrt{\frac{\log(pq)}{n}}; \quad \mathbb{R}_0 := \max_{k \in \mathcal{I}_K} \mathbb{R}(v_\Theta, \Sigma_x^k, \Sigma_y^k)$$

the following holds

$$\|\hat{\beta} - \beta_0\|_1 \leq 48\sqrt{|h_{\max}|} s_\beta \lambda_n / \psi^* \quad (\text{A.2})$$

$$\|\hat{\beta} - \beta_0\| \leq 12\sqrt{s_\beta} \lambda_n / \psi^* \quad (\text{A.3})$$

$$\sum_{h \in \mathcal{H}} \|\beta^{[h]} - \beta_0^{[h]}\| \leq 48s_\beta \lambda_n / \psi^* \quad (\text{A.4})$$

$$(\hat{\beta} - \beta_0)^T \hat{\Gamma} (\hat{\beta} - \beta_0) \leq 72s_\beta \lambda_n^2 / \psi^* \quad (\text{A.5})$$

Next we verify that conditions (B2) and (B3) hold with high probability given fixed  $\hat{\Theta}$ .

**Proposition A.5.** Consider deterministic  $\hat{\Theta}$  satisfying assumption (B1). Assume that the matrices  $(\hat{\mathbf{T}}^k)^2, k \in \mathcal{I}_K$  are diagonally dominant. Then for sample size  $n \gtrsim \log(pq)$ ,

1.  $\hat{\Gamma}$  satisfies the RE condition:  $\hat{\Gamma} \sim RE(\psi_*, \phi_*)$ , where

$$\psi_* = \min_k \psi^k \left( \min_i \psi_t^i - dv_\Theta \right), \phi_* = \max_k \phi^k \left( \min_i \phi_t^i + dv_\Theta \right)$$

with probability  $\geq 1 - 2 \exp(-c_3 n)$ ,  $c_3 > 0$ .

2. The deviation bound in (B2) is satisfied with probability  $\geq 1 - 12c_1 \exp[-(c_2^2 - 1) \log(pq)]$ ,  $c_1 > 0, c_2 > 1$ , where

$$\mathbb{R}(v_\Theta, \Sigma_x^k, \Sigma_y^k) = c_2 \sqrt{\Lambda_{\max}(\Sigma_x^k)} \left( dv_\Theta \Lambda_{\min}(\Sigma_y^k) + \frac{1}{\Lambda_{\min}(\Sigma_y^k)} \right)$$

We now put both the pieces together, and prove that our alternating algorithm results in a solution sequence  $\{\hat{\beta}^{(r)}, \hat{\Theta}^{(r)}\}, r = 1, 2, \dots$  that lies uniformly within a non-expanding ball around the true parameter values.

*Proof of Theorem A.1.* The proof has three parts, where we prove the consistency of the neighborhood regression coefficients, selection of edge sets, and finally the refitting step, respectively. This is the same structure as the proof of Theorem 1 in [Ma and Michailidis \(2016\)](#), where they prove consistency of the (single layer) JSEM estimates. The derivation of the first part is different from that in the JSEM proof, which we shall show in detail (in the proof of Proposition A.6). The second and third parts follow similar lines, incorporating the updated quantities from the first part. For these we provide outlines and leave the details to the reader.

*Step 1: consistency of neighborhood regression.* The following proposition establishes error bounds for estimated neighborhood coefficients in the Y-network.

**Proposition A.6.** *Consider the estimation problem in (2.13) and take  $\gamma_n \geq 4\sqrt{|g_{\max}|}\mathbb{Q}_0$ . Given the conditions (T2) and (T3) hold, for any solution of (2.13) we shall have*

$$\|\hat{\Theta}_j - \Theta_{0,j}\|_F \leq 12\sqrt{s_j}\gamma_n/\psi \quad (\text{A.6})$$

$$\sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\hat{\boldsymbol{\theta}}_{jj'}^{[g]} - \boldsymbol{\theta}_{0,jj'}^{[g]}\| \leq 48s_j\gamma_n/\psi \quad (\text{A.7})$$

Also denote the non-zero support of  $\hat{\Theta}_j$  by  $\hat{\mathcal{S}}_j$ , i.e.  $\hat{\mathcal{S}}_j = \{(j', g) : \hat{\boldsymbol{\theta}}_{jj'}^{[g]} \neq \mathbf{0}\}$ . Then

$$|\hat{\mathcal{S}}_j| \leq 128s_j/\psi \quad (\text{A.8})$$

*Step 2: Edge set selection.* We denote the selected edge set for the  $k^{\text{th}}$  Y-network by  $\hat{E}^k$ . Denote its population version by  $E_0^k$ . Further, let

$$\tilde{\Omega}_y^k = \text{diag}(\Omega_y^k) + \Omega_{y, E_0^k \cap \hat{E}^k}^k$$

With similar derivations to the proof of Corollary A.1 in [Ma and Michailidis \(2016\)](#), The following two upper bounds can be established:

$$|\hat{E}^k| \leq \frac{128S}{\psi} \quad (\text{A.9})$$

$$\frac{1}{K} \sum_{k=1}^K \|\tilde{\Omega}_y^k - \Omega_y^k\|_F \leq \frac{12c_0\sqrt{S}\gamma_n}{\sqrt{K}\psi} \quad (\text{A.10})$$

following which, taking  $\gamma_n = 4\sqrt{|g_{\max}|}\mathbb{Q}_0$ ,

$$\Lambda_{\min}(\tilde{\Omega}_y^k) \geq d_0 - \frac{48c_0\mathbb{Q}_0\sqrt{|g_{\max}|}S}{\psi} \geq (1 - t_1)d_0 > 0 \quad (\text{A.11})$$

$$\Lambda_{\max}(\tilde{\Omega}_y^k) \leq c_0 + \frac{48c_0\mathbb{Q}_0\sqrt{|g_{\max}|}S}{\psi} \leq c_0 + t_1d_0 < \infty \quad (\text{A.12})$$

with  $0 < t_1 < 1$ , and the sample size  $n$  satisfying

$$n \geq |g_{\max}|S \left[ \frac{48c_0\mathbb{Q}_0}{\psi t_1 d_0} \right]^2; \quad \mathbb{Q}_0 := \sqrt{n}\mathbb{Q}_0$$

*Step 3: Refitting.* Following the same steps as part A.3 in the proof of Theorem 4.1 in [Ma and Michailidis \(2016\)](#), it can be proven using (A.9)–(A.12) that

$$\sum_{k=1}^K \|\hat{\Omega}_y^k - \tilde{\Omega}_y^k\|^2 \leq O(\mathbb{Q}_0^2 |g_{\max}| S)$$

The proof is now complete by combining this with (A.10) then applying Cauchy-Schwarz inequality and triangle inequality.  $\square$

*Proof of Proposition A.2.* We drop the superscript  $k$  since there is no scope of ambiguity. For part 1, we start with an auxiliary lemma:

**Lemma A.7.** *For a sub-gaussian design matrix  $\mathbf{X} \in \mathbb{M}(n, p)$  with columns having mean  $\mathbf{0}_p$  and covariance matrix  $\Sigma_x$ , the sample covariance matrix  $\hat{\Sigma}_x = \mathbf{X}^T \mathbf{X} / n$  satisfies the RE condition*

$$\hat{\Sigma}_x \sim RE \left( \frac{\Lambda_{\min}(\Sigma_x)}{2}, \frac{\Lambda_{\min}(\Sigma_x) \log p}{2n} \right)$$

with probability  $\geq 1 - 2 \exp(-c_3 n)$  for some  $c_3 > 0$ .

Now denote  $\hat{\mathbf{E}} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}$ . For  $\mathbf{v} \in \mathbb{R}^q$ , we have

$$\begin{aligned} \mathbf{v}^T \hat{\mathbf{S}} \mathbf{v} &= \frac{1}{n} \|\hat{\mathbf{E}} \mathbf{v}\|^2 \\ &= \frac{1}{n} \|(\mathbf{E} + \mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}})) \mathbf{v}\|^2 \\ &= \mathbf{v}^T \mathbf{S} \mathbf{v} + \frac{1}{n} \|\mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}}) \mathbf{v}\|^2 + 2 \mathbf{v}^T (\mathbf{B}_0 - \hat{\mathbf{B}})^T \left( \frac{(\mathbf{X})^T \mathbf{E}}{n} \right) \mathbf{v} \end{aligned} \quad (\text{A.13})$$

For the first summand,  $\mathbf{v}^T \mathbf{S}^k \mathbf{v} \geq \psi_y \|\mathbf{v}\|^2 - \phi_y \|\mathbf{v}\|_1^2$  with  $\psi_y = \Lambda_{\min}(\Sigma_y)/2$ ,  $\phi_y = \psi_y \log p/n$  by applying Lemma A.7 on  $\mathbf{S}$ . The second summand is greater than or equal to 0. For the third summand,

$$2 \mathbf{v}^T (\mathbf{B}_0 - \hat{\mathbf{B}})^T \left( \frac{(\mathbf{X})^T \mathbf{E}}{n} \right) \mathbf{v} \geq -2v_\beta \left\| \frac{(\mathbf{X})^T \mathbf{E}}{n} \right\|_\infty \|\mathbf{v}\|_1^2$$

by assumption (T1). Now we use another lemma:

**Lemma A.8.** *For zero-mean independent sub-gaussian matrices  $\mathbf{X} \in \mathbb{M}(n, p)$ ,  $\mathbf{E} \in \mathbb{M}(n, q)$  with parameters  $(\Sigma_x, \sigma_x^2)$  and  $(\Sigma_e, \sigma_e^2)$  respectively, given that  $n \gtrsim \log(pq)$  the following holds with probability  $\geq 1 - 6c_1 \exp[-(c_2^2 - 1) \log(pq)]$  for some  $c_1 > 0$ ,  $c_2 > 1$ :*

$$\frac{1}{n} \|\mathbf{X}^T \mathbf{E}\|_\infty \leq c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_e)]^{1/2} \sqrt{\frac{\log(pq)}{n}}$$

Subsequently we collect all summands in (A.13) and get

$$\mathbf{v}^T \hat{\mathbf{S}} \mathbf{v} \geq \psi_y \|\mathbf{v}\|^2 - \left( \phi_y + 2v_\beta c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_y)]^{1/2} \sqrt{\frac{\log(pq)}{n}} \right) \|\mathbf{v}\|_1^2$$

with probability  $\geq 1 - 2 \exp(-c_3 n) - 6c_1 \exp[-(c_2^2 - 1) \log(pq)]$ . This concludes the proof of part 1.

To prove part 2, we decompose the quantity in question:

$$\begin{aligned}
\left\| \frac{1}{n} \widehat{\mathbf{E}}_{-j}^T \widehat{\mathbf{E}} \mathbf{T}_{0,j} \right\|_{\infty} &= \left\| \frac{1}{n} \left[ \mathbf{E}_{-j} + \mathbf{X}(\mathbf{B}_{0,j} - \widehat{\mathbf{B}}_j) \right]^T \left[ \mathbf{E} + \mathbf{X}(\mathbf{B}_0 - \widehat{\mathbf{B}}) \right] \mathbf{T}_{0,j} \right\|_{\infty} \\
&\leq \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{E} \mathbf{T}_{0,j} \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{X}(\mathbf{B}_0 - \widehat{\mathbf{B}}) \mathbf{T}_{0,j} \right\|_{\infty} \\
&\quad + \left\| \frac{1}{n} (\mathbf{B}_{0,j} - \widehat{\mathbf{B}}_j)^T \mathbf{X}^T \mathbf{X}(\mathbf{B}_0 - \widehat{\mathbf{B}}) \mathbf{T}_{0,j} \right\|_{\infty} + \left\| \frac{1}{n} (\mathbf{B}_{0,j} - \widehat{\mathbf{B}}_j)^T \mathbf{X}^T \mathbf{E} \mathbf{T}_{0,j} \right\|_{\infty} \\
&= \|\mathbf{W}_1\|_{\infty} + \|\mathbf{W}_2\|_{\infty} + \|\mathbf{W}_3\|_{\infty} + \|\mathbf{W}_4\|_{\infty} \tag{A.14}
\end{aligned}$$

Now

$$\mathbf{W}_1 = \frac{1}{n} \mathbf{E}_{-j}^T (\mathbf{E}_j - \mathbf{E}_{-j} \boldsymbol{\theta}_{0,j})$$

For node  $j$  in the  $y$ -network,  $\mathbb{E}_{-j}$  and  $E_j - \mathbb{E}_{-j} \boldsymbol{\theta}_{0,j}$  are the neighborhood regression coefficients and residuals, respectively. Thus they are orthogonal, so we can apply Lemma A.8 on  $\mathbf{E}_{-j}$  and  $\mathbf{E}_j - \mathbf{E}_{-j} \boldsymbol{\theta}_{0,j}$  to obtain that for  $n \gtrsim \log(q-1)$ ,

$$\|\mathbf{W}_1\|_{\infty} \leq c_5 [\Lambda_{\max}(\Sigma_{y,-j}) \sigma_{y,j,-j}]^{1/2} \sqrt{\frac{\log(q-1)}{n}} \tag{A.15}$$

holds with probability  $\geq 1 - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$  for some  $c_4 > 0, c_5 > 1$ .

The same bounds hold for  $\mathbf{W}_2$  and  $\mathbf{W}_4$ :

$$\begin{aligned}
\|\mathbf{W}_2\|_{\infty} &\leq \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{X}(\mathbf{B}_0 - \widehat{\mathbf{B}}) \right\|_{\infty} \|\mathbf{T}_{0,j}\|_1 \leq \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{X} \right\|_{\infty} \|\mathbf{B}_0 - \widehat{\mathbf{B}}\|_1 \|\mathbf{T}_{0,j}\|_1 \\
\|\mathbf{W}_4\|_{\infty} &\leq \left\| \frac{1}{n} (\mathbf{B}_{0,j} - \widehat{\mathbf{B}}_j)^T \mathbf{X}^T \mathbf{E} \right\|_{\infty} \|\mathbf{T}_{0,j}\|_1 \leq \left\| \frac{1}{n} \mathbf{E}^T \mathbf{X} \right\|_{\infty} \|\mathbf{B}_0 - \widehat{\mathbf{B}}\|_1 \|\mathbf{T}_{0,j}\|_1
\end{aligned}$$

Now since  $\Omega_y$  is diagonally dominant,  $|\omega_{y,jj}| \geq \sum_{j \neq j'} |\omega_{y,jj'}|$  for any  $j \in \mathcal{I}_q$ . Hence

$$\|\mathbf{T}_{0,j}\|_1 = \sum_{j'=1}^q |T_{jj'}| = 1 + \sum_{j \neq j'} |\theta_{jj'}| = 1 + \frac{1}{\omega_{y,jj}} \sum_{j \neq j'} |\omega_{y,jj'}| \leq 2$$

so that for  $n \gtrsim \log(pq)$ ,

$$\|\mathbf{W}_2\|_{\infty} + \|\mathbf{W}_4\|_{\infty} \leq 4v_{\beta} c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_y)]^{1/2} \sqrt{\frac{\log(pq)}{n}} \tag{A.16}$$

with probability  $\geq 1 - 6c_1 \exp[-(c_2^2 - 1) \log(pq)]$  by applying Lemma A.8 and assumption (T1).

Finally for  $\mathbf{W}_3$ , we apply Lemma 8 of Ravikumar et al. (2011) on the (sub-gaussian) design matrix  $\mathbf{X}$  to obtain that for sample size

$$n \geq 512(1 + 4\Lambda_{\max}(\Sigma_x^k))^4 \max_i (\sigma_{x,ii}^k)^4 \log(4p^{\tau_1}) \tag{A.17}$$

we get that with probability  $\geq 1 - 1/p^{\tau_1-2}$ ,  $\tau_1 > 2$ ,

$$\left\| \frac{\mathbf{X}^T \mathbf{X}}{n} \right\|_{\infty} \leq \sqrt{\frac{\log 4 + \tau_1 \log p}{c_x n}} + \max_i \sigma_{x,ii} = V_x; \quad c_x = \left[ 128(1 + 4\Lambda_{\max}(\Sigma_x))^2 \max_i (\sigma_{x,ii})^2 \right]^{-1}$$

Thus with the same probability,

$$\|\mathbf{W}_4\|_{\infty} \leq \left\| \frac{\mathbf{X}^T \mathbf{X}}{n} \right\|_{\infty} \|\widehat{\mathbf{B}} - \mathbf{B}_0\|_1^2 \|\mathbf{T}_{0,j}\|_1 \leq 2v_{\beta}^2 V_x \quad (\text{A.18})$$

We now bound the right hand side of (A.14) using (A.15), (A.16) and (A.18) to complete the proof, with the leading term of the sample size requirement being  $n \gtrsim \log(pq)$ .  $\square$

*Proof of Theorem A.4.* The proof follows that of Proposition A.6, with a different group norm structure. We only point out the differences.

Putting  $\beta = \beta_0$  in (2.12) we get

$$-2\widehat{\beta}^T \widehat{\gamma} + \beta^T \widehat{\Gamma} \widehat{\beta} + \lambda_n \sum_{h \in \mathcal{H}} \|\widehat{\beta}^{[h]}\| \leq -2\beta_0^T \widehat{\gamma} + \beta_0^T \widehat{\Gamma} \beta_0 + \lambda_n \sum_{h \in \mathcal{H}} \|\beta_0^{[h]}\|$$

Denote  $\mathbf{b} = \widehat{\beta} - \beta_0$ . Then we have

$$\mathbf{b}^T \widehat{\Gamma} \mathbf{b} \leq 2\mathbf{b}^T (\widehat{\gamma} - \widehat{\Gamma} \beta_0) + \lambda_n \sum_{h \in \mathcal{H}} (\|\beta_0^{[h]}\| - \|\beta_0^{[h]}\| + \mathbf{b}^{[h]})$$

Proceeding similarly as the proof of Proposition A.6, with a different deviation bound and choice of  $\lambda_n$ , we get expressions equivalent to (B.3) and (B.4) respectively:

$$\mathbf{b}^T \widehat{\Gamma} \mathbf{b} \leq \frac{3}{2} \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\| \quad (\text{A.19})$$

$$\frac{\psi^*}{3} \|\mathbf{b}\|^2 \leq \lambda_n \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\| \leq 4\lambda_n \sqrt{s_{\beta}} \|\mathbf{b}\| \quad (\text{A.20})$$

Furthermore,  $\|\mathbf{b}\|_1 \leq \sqrt{|h_{\max}|} \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\|$ . The bounds in (A.2), (A.3), (A.4) and (A.5) now follow.  $\square$

*Proof of Proposition A.5.* For part 1 it is enough to prove that with  $\widehat{\Sigma}_x^k := (\mathbf{X}^k)^T \mathbf{X}^k / n$ ,

$$\widehat{\mathbf{T}}_k^2 \otimes \widehat{\Sigma}_x^k \sim RE(\psi_*^k, \phi_*^k) \quad (\text{A.21})$$

with high enough probability. because then we can take  $\psi_* = \min_k \psi_*^k$ ,  $\phi_* = \max_k \phi_*^k$ . The proof of (A.21) follows similar lines of the proof of Proposition 1 in Lin et al. (2016), only replacing  $\Theta_{\epsilon}$ ,  $\widehat{\Theta}_{\epsilon}$ ,  $\mathbf{X}$  therein with  $(\mathbf{T}^k)^2$ ,  $(\widehat{\mathbf{T}}^k)^2$ ,  $\mathbf{X}^k$ , respectively. We omit the details.

Part 2 follows the proof of Proposition 2 in Lin et al. (2016).  $\square$

We also have the following auxiliary result that helps establish results in the testing section.

**Lemma A.9.** lemma and proof tbd

*Proof of Theorem 3.1.* Let us define the following:

$$\begin{aligned}\widehat{\Omega}_y &= \text{diag}(\widehat{\Omega}_y^1, \widehat{\Omega}_y^2) \\ \mathbf{M}_i &= \text{diag}(m_i^1, m_i^2) \\ \widehat{\mathbf{C}}_i &= \text{diag}(\widehat{\mathbf{c}}_i^1, \widehat{\mathbf{c}}_i^2) \\ \widehat{\mathbf{D}}_i &= \text{diag}(\widehat{\mathbf{b}}_i^1, \widehat{\mathbf{b}}_i^2) \\ \mathbf{D}_i &= \text{diag}(\mathbf{b}_{0,i}^1, \mathbf{b}_{0,i}^2) \\ \mathbf{R}_i^k &= \mathbf{X}_i^k - \mathbf{X}_{-i}^k \widehat{\boldsymbol{\zeta}}_i^k; k = 1, 2\end{aligned}$$

Then from (3.1) we have

$$\mathbf{M}_i(\widehat{\mathbf{C}}_i - \widehat{\mathbf{D}}_i)^T = \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \widehat{\mathbf{E}}^1 \\ \frac{1}{\widehat{s}_i^2} (\mathbf{R}_i^2)^T \widehat{\mathbf{E}}^2 \end{bmatrix} \quad (\text{A.22})$$

We now decompose  $\widehat{\mathbf{E}}^k$ :

$$\begin{aligned}\widehat{\mathbf{E}}^k &= \mathbf{Y}^k - \mathbf{X}^k \widehat{\mathbf{B}}^k \\ &= \mathbf{E}^k + \mathbf{X}^k (\mathbf{B}_0^k - \widehat{\mathbf{B}}^k) \\ &= \mathbf{E}^k + \mathbf{X}_i^k (\mathbf{b}_{0,i}^k - \widehat{\mathbf{b}}_i^k) + \mathbf{X}_{-i}^k (\mathbf{B}_{0,-i}^k - \widehat{\mathbf{B}}_{-i}^k)\end{aligned}$$

Putting them back in (A.22) and using  $t^k = (\mathbf{R}^k)^T \mathbf{X}^k / n$ ,

$$\begin{aligned}\mathbf{M}_i(\widehat{\mathbf{C}}_i - \widehat{\mathbf{D}}_i)^T &= \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{E}^1 \\ \frac{1}{\widehat{s}_i^2} (\mathbf{R}_i^2)^T \mathbf{E}^2 \end{bmatrix} + \mathbf{M}_i(\mathbf{D}_i - \widehat{\mathbf{D}}_i)^T + \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{X}_{-i}^1 (\mathbf{B}_{0,-i}^1 - \widehat{\mathbf{B}}_{-i}^1) \\ \frac{1}{\widehat{s}_i^2} (\mathbf{R}_i^2)^T \mathbf{X}_{-i}^2 (\mathbf{B}_{0,-i}^2 - \widehat{\mathbf{B}}_{-i}^2) \end{bmatrix} \\ \Rightarrow \widehat{\Omega}_y^{1/2} \mathbf{M}_i(\widehat{\mathbf{C}}_i - \mathbf{D}_i)^T &= \frac{\widehat{\Omega}_y^{1/2}}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{E}^1 \\ \frac{1}{\widehat{s}_i^2} (\mathbf{R}_i^2)^T \mathbf{E}^2 \end{bmatrix} + \frac{\widehat{\Omega}_y^{1/2}}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{X}_{-i}^1 (\mathbf{B}_{0,-i}^1 - \widehat{\mathbf{B}}_{-i}^1) \\ \frac{1}{\widehat{s}_i^2} (\mathbf{R}_i^2)^T \mathbf{X}_{-i}^2 (\mathbf{B}_{0,-i}^2 - \widehat{\mathbf{B}}_{-i}^2) \end{bmatrix}\end{aligned} \quad (\text{A.23})$$

At this point, we drop  $k$  in the subscripts, and prove the following:

**Lemma A.10.** Given conditions (C1) and (C2), the following holds for sample size  $n$  such that  $n \gtrsim \log(pq)$  and  $\sigma_{x,i,-i} - n^{-1/4} - v_\zeta \sqrt{V_x} > 0$ :

$$\begin{aligned}\frac{1}{\sqrt{n \widehat{s}_i}} \widehat{\Omega}_y^{1/2} \mathbf{E}^T \mathbf{R}_i &\sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{1n}; \\ \|\mathbf{S}_{1n}\|_\infty &\leq \frac{v_\Omega(2 + v_\zeta) c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_e)]^{1/2} \sqrt{\log(pq)}}{\sigma_{x,i,-i} - n^{-1/4} - v_\zeta \sqrt{V_x}} = O\left(\frac{\log(pq)}{\sqrt{n}}\right)\end{aligned} \quad (\text{A.24})$$



with probability larger than or equal to

$$1 - 6c_1 e^{-(c_2^2 - 1) \log pq} - \frac{1}{p^{\tau_1 - 2}} - \frac{\kappa_i}{\sqrt{n}} \quad (\text{A.25})$$

for some  $c_1, c_4 > 0$ ,  $c_2, c_5 > 1$ , and  $\kappa_i := \mathbb{V}[(X_i - \mathbb{X}_{-i} \boldsymbol{\zeta}_{0,-i})^2]$ . Additionally, given condition (C3)

$$\begin{aligned} & \left\| \frac{1}{\sqrt{n} \hat{s}_i} \mathbf{R}_i^T \mathbf{X}_{-i} (\mathbf{B}_{0,-i} - \hat{\mathbf{B}}_{-i}) \hat{\Omega}_y^{1/2} \right\|_{\infty} \\ & \leq \frac{v_{\beta} (\Lambda_{\min}(\Sigma_y)^{1/2} + v_{\Omega})}{\sigma_{x,i,-i} - n^{-1/2} - v_{\zeta} \sqrt{V_x}} \left[ c_7 \sqrt{(\sigma_{x,i,-i} \Lambda_{\max}(\Sigma_{x,-i})) \log p} + \sqrt{n} v_{\zeta} V_x \right] = O\left(\frac{\log(pq)}{\sqrt{n}}\right) \end{aligned} \quad (\text{A.26})$$

with probability condition (A.25).

Given Lemma A.10, the first and second summands on the right hand side of (A.23) are bounded above by applying each of (A.24) and (A.26) twice, respectively. This completes our proof.  $\square$

*Proof of Theorem 3.4.*  $\square$

## B Proofs of auxiliary results

### need modification

*Proof of Proposition A.6.* In its reparametrized version, (2.13) becomes

$$\hat{\mathbf{T}}_j = \arg \min_{\mathbf{T}_j} \left\{ \frac{1}{n} \sum_{k=1}^K \|(\mathbf{Y}^k - \mathbf{X}^k \hat{\mathbf{B}}^k) \mathbf{T}_j^k\|^2 + \gamma_n \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{T}_{jj'}^{[g]}\| \right\} \quad (\text{B.1})$$

with  $\mathbf{T}_{jj'}^{[g]} := (T_{jj'}^k)_{k \in g}$ . Now for any  $\mathbf{T}_j \in \mathbb{M}(q, K)$  we have

$$\frac{1}{n} \sum_{k=1}^K \|(\mathbf{Y}^k - \mathbf{X}^k \hat{\mathbf{B}}^k) \hat{\mathbf{T}}_j^k\|^2 + \gamma_n \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\hat{\mathbf{T}}_{jj'}^{[g]}\| \leq \frac{1}{n} \sum_{k=1}^K \|(\mathbf{Y}^k - \mathbf{X}^k \hat{\mathbf{B}}^k) \mathbf{T}_j^k\|^2 + \gamma_n \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{T}_{jj'}^{[g]}\|$$

For  $\mathbf{T}_j = \mathbf{T}_{0,j}$  this reduces to

$$\sum_{k=1}^K (\mathbf{d}_j^k)^T \hat{\mathbf{S}}^k \mathbf{d}_j^k \leq -2 \sum_{k=1}^K (\mathbf{d}_j^k)^T \hat{\mathbf{S}}^k \mathbf{T}_{0,j}^k + \gamma_n \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \left( \|\mathbf{T}_{jj'}^{[g]}\| - \|\mathbf{T}_{jj'}^{[g]} + \mathbf{d}_{jj'}^{[g]}\| \right) \quad (\text{B.2})$$

with  $\mathbf{d}_j^k := \widehat{\mathbf{T}}_j^k - \mathbf{T}_{0,j}^k$  etc. For the  $k^{\text{th}}$  summand in the first term on the right hand side, since  $d_{jj}^k = 0$ ,  $\widehat{\mathbf{E}}^k \mathbf{d}_j^k = \widehat{\mathbf{E}}_{-j}^k \mathbf{d}_{-j}^k$ . Thus

$$\begin{aligned} \sum_{k=1}^K \left| (\mathbf{d}_j^k)^T \widehat{\mathbf{S}}^k \mathbf{T}_{0,j}^k \right| &= \sum_{k=1}^K \left| \mathbf{d}_j^k \cdot \frac{1}{n} (\widehat{\mathbf{E}}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right| \\ &\leq \sum_{k=1}^K \left\| \frac{1}{n} (\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_{\infty} \|\mathbf{d}_{-j}^k\|_1 \\ &\leq Q_0 \sqrt{|g_{\max}|} \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \end{aligned}$$

by assumption (T2). For the second term, suppose  $\mathcal{S}_{0,j}$  is the support of  $\Theta_{0,j}$ , i.e.  $\mathcal{S}_{0,j} = \{(j', g) : \theta_{jj'}^{[g]} \neq 0\}$ . Then

$$\begin{aligned} \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \left( \|\mathbf{T}_{jj'}^{[g]}\| - \|\mathbf{T}_{jj'}^{[g]} + \mathbf{d}_{jj'}^{[g]}\| \right) &\leq \sum_{(j', g) \in \mathcal{S}_{0,j}} \left( \|\mathbf{T}_{jj'}^{[g]}\| - \|\mathbf{T}_{jj'}^{[g]} + \mathbf{d}_{jj'}^{[g]}\| \right) - \sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \\ &\leq \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| - \sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \end{aligned}$$

so that by choice of  $\gamma_n$  (B.2) reduces to

$$\begin{aligned} \sum_{k=1}^K (\mathbf{d}_j^k)^T \widehat{\mathbf{S}}^k \mathbf{d}_j^k &\leq \frac{\gamma_n}{2} \left[ \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| + \sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \right] + \gamma_n \left[ \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| - \sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \right] \\ &= \frac{3\gamma_n}{2} \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| - \frac{\gamma_n}{2} \sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \\ &\leq \frac{3\gamma_n}{2} \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \end{aligned} \tag{B.3}$$

Since the left hand side is  $\geq 0$ , this also implies

$$\sum_{(j', g) \notin \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \leq 3 \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \Rightarrow \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \leq 4 \sum_{(j', g) \in \mathcal{S}_{0,j}} \|\mathbf{d}_{jj'}^{[g]}\| \leq 4\sqrt{s_j} \|\mathbf{D}_j\|_F$$

with  $\mathbf{D}_j = (\mathbf{d}_j^k)_{k \in \mathcal{I}_K}$ . Now the RE condition on  $\widehat{\mathbf{S}}^k$  means that

$$\sum_{k=1}^K (\mathbf{d}_j^k)^T \widehat{\mathbf{S}}^k \mathbf{d}_j^k \geq \sum_{k=1}^K \left( \psi_k \|\mathbf{d}_j^k\|^2 - \phi_k \|\mathbf{d}_j^k\|_1^2 \right) \geq \psi \|\mathbf{D}_j\|_F^2 - \phi \|\mathbf{D}_j\|_1^2 \geq (\psi - Kq\phi) \|\mathbf{D}_j\|_F^2 \geq \frac{\psi}{2} \|\mathbf{D}_j\|_F^2$$

by assumption (T3).

Thus we finally have

$$\frac{\psi}{3} \|\mathbf{D}_j\|_F^2 \leq \gamma_n \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \leq 4\gamma_n \sqrt{s_j} \|\mathbf{D}_j\|_F \tag{B.4}$$

Since

$$(\mathbf{D}_j)_{j',k} = \widehat{T}_{jj'}^k - T_{0,jj'}^k = \begin{cases} 0 & \text{if } j = j' \\ -(\widehat{\theta}_{jj'}^k - \theta_{0,jj'}^k) & \text{if } j \neq j' \end{cases}$$

The bounds in (A.6) and (A.7) are obtained by replacing the corresponding elements in (B.4).

For the bound on  $|\widehat{\mathcal{S}}_j|$ , notice that if  $\widehat{\theta}_{jj'}^{[g]} \neq 0$  for some  $(j', g)$ ,

$$\begin{aligned} \frac{1}{n} \sum_{k \in g} \left| ((\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k (\widehat{\mathbf{T}}_j^k - \mathbf{T}_{0,j}^k))^{j'} \right| &\geq \frac{1}{n} \sum_{k \in g} \left| ((\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \widehat{\mathbf{T}}_j^k)^{j'} \right| - \frac{1}{n} \sum_{k \in g} \left| ((\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k)^{j'} \right| \\ &\geq |g| \gamma_n - \sum_{k \in g} \mathbb{Q}(v_\beta, \Sigma_x^k, \Sigma_y^k) \end{aligned}$$

using the KKT condition for (2.13) and assumption (T2). The choice of  $\gamma_n$  now ensures that the right hand side is  $\geq 3|g|\gamma_n/4$ . Hence

$$\begin{aligned} |\widehat{\mathcal{S}}_j| &\leq \sum_{(j',g) \in \widehat{\mathcal{S}}_j} \frac{16}{9n^2|g|^2\gamma_n^2} \sum_{k \in g} \left| ((\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k (\widehat{\mathbf{T}}_j^k - \mathbf{T}_{0,j}^k))^{j'} \right|^2 \\ &\leq \frac{16}{9\gamma_n^2} \sum_{k=1}^K \frac{1}{n} \left\| (\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k (\widehat{\mathbf{T}}_j^k - \mathbf{T}_{0,j}^k) \right\|^2 \\ &= \frac{16}{9\gamma_n^2} \sum_{k=1}^K (\mathbf{d}_j^k)^T \widehat{\mathbf{S}}^k \mathbf{d}_j^k \\ &\leq \frac{8}{3\gamma_n} \sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \leq \frac{128s_j}{\psi} \end{aligned}$$

using (B.3) and (B.4).  $\square$

*Proof of Lemma A.7.* This is same as Lemma 2 in Appendix B of Lin et al. (2016) and its proof can be found there.  $\square$

*Proof of Lemma A.8.* This is a part of Lemma 3 of Appendix B in Lin et al. (2016), and is proved therein.  $\square$

*Proof of Lemma A.10.* To show (A.24) we have

$$\frac{1}{\sqrt{n\widehat{s}_i}} \widehat{\Omega}_y^{1/2} \mathbf{E}^T \mathbf{R}_i = \frac{1}{\sqrt{n\widehat{s}_i}} (\widehat{\Omega}_y^{1/2} - \Omega_y^{1/2}) \mathbf{E}^T \mathbf{R}_i + \frac{1}{\sqrt{n\widehat{s}_i}} \Omega_y^{1/2} \mathbf{E}^T \mathbf{R}_i$$

The second summand is distributed as  $\mathcal{N}_q(\mathbf{0}, \mathbf{I})$ . For the first summand,

$$\begin{aligned}
\frac{1}{\sqrt{n}} \left\| (\hat{\Omega}_y^{1/2} - \Omega_y^{1/2}) \mathbf{E}^T \mathbf{R}_i \right\|_\infty &\leq \frac{1}{\sqrt{n}} \left\| \hat{\Omega}_y^{1/2} - \Omega_y^{1/2} \right\|_\infty \left\| \mathbf{E}^T \mathbf{R}_i \right\|_1 \\
&\leq \sqrt{n} v_\Omega \frac{1}{n} \left[ \left\| \mathbf{E}^T (\mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_i) \right\|_1 + \left\| \mathbf{E}^T \mathbf{X}_{-i} (\hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i}) \right\|_1 \right] \\
&\leq \sqrt{n} v_\Omega \frac{1}{n} \left[ \left\| \mathbf{E}^T \mathbf{X}_i \right\|_\infty + \left\| \mathbf{E}^T \mathbf{X}_{-i} \right\|_\infty \left\{ \left\| \boldsymbol{\zeta}_i \right\|_1 + \left\| \hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_i \right\|_1 \right\} \right] \\
&\leq \sqrt{n} v_\Omega \left[ \frac{1}{n} \left\| \mathbf{E}^T \mathbf{X}_i \right\|_\infty + \frac{1 + v_\zeta}{n} \left\| \mathbf{E}^T \mathbf{X}_{-i} \right\|_\infty \right] \\
&\leq \sqrt{n} v_\Omega (2 + v_\zeta) \cdot \frac{1}{n} \left\| \mathbf{E}^T \mathbf{X} \right\|_\infty
\end{aligned}$$

because  $\Omega_x$  is diagonally dominant implies  $\left\| \boldsymbol{\zeta}_i \right\|_1 = \sum_{i' \neq i} |\omega_{x,ii'}| / \omega_{x,ii} \leq 1$ , and using assumption (C1). Applying Lemma A.8, the following holds:

$$\frac{1}{\sqrt{n}} \left\| (\hat{\Omega}_y^{1/2} - \Omega_y^{1/2}) \mathbf{E}^T \mathbf{R}_i \right\|_\infty \leq v_\Omega (2 + v_\zeta) c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_e)]^{1/2} \sqrt{\log(pq)} \quad (\text{B.5})$$

with probability  $\geq 1 - 6c_1 \exp[-(c_2^2 - 1) \log(pq)]$  for some  $c_1 > 0, c_2 > 1$ .

On the other hand

$$s_i^2 := \frac{1}{n} \left\| \mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_i \right\|^2 \leq \hat{s}_i^2 + \frac{1}{n} \left\| \mathbf{X}_{-i} (\hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i}) \right\|^2 \leq \hat{s}_i^2 + \left\| \hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i} \right\|_1^2 \left\| \frac{1}{n} \mathbf{X}_{-i}^T \mathbf{X}_{-i} \right\|_\infty$$

which implies  $s_i \leq \hat{s}_i + v_\zeta \sqrt{V_x}$ . By applying Lemma 8 of Ravikumar et al. (2011),

$$\left\| \frac{1}{n} \mathbf{X}_{-i}^T \mathbf{X}_{-i} \right\|_\infty \leq \left\| \frac{1}{n} \mathbf{X}^T \mathbf{X} \right\|_\infty \leq V_x \quad (\text{B.6})$$

with probability  $\geq 1 - 1/p^{\tau_1 - 2}, \tau_1 > 2$ , and

$$n \geq 512(1 + 4\Lambda_{\max}(\Sigma_x))^4 \max_i (\sigma_{x,ii})^4 \log(4p^{\tau_1}) \quad (\text{B.7})$$

On the other hand, by Chebyshev inequality, for any  $\epsilon > 0$

$$P(|s_i - \sigma_{x,i,-i}| \geq \epsilon) \leq \frac{\mathbb{V} s_i}{\epsilon^2} = \frac{\kappa_i}{n\epsilon^2}$$

Taking  $\epsilon = n^{-1/4}$ , we have  $s_i \geq \sigma_{x,i,-i} - n^{-1/4}$  with probability  $\geq 1 - \kappa_i n^{-1/2}$ . Then, for  $n$  satisfying (B.6) and  $\sigma_{x,i,-i} - n^{-1/4} > v_\zeta \sqrt{V_x}$ , we get the bound with the above probability:

$$\frac{1}{\hat{s}_i} \leq \frac{1}{\sigma_{x,i,-i} - n^{-1/4} - v_\zeta \sqrt{V_x}} \quad (\text{B.8})$$

Combining (B.5) and (B.8) gives the upper bound for the right hand side of (A.24) with the requisite probability and sample size conditions.

To prove (A.26) we have

$$\frac{1}{n} \left\| \mathbf{R}_i^T \mathbf{X}_{-i} \right\|_\infty \leq \frac{1}{n} \left\| (\mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_{0,i})^T \mathbf{X}_{-i} \right\|_\infty + \frac{1}{n} \left\| \mathbf{X}_{-i}^T \mathbf{X}_{-i} (\hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i}) \right\|_\infty \quad (\text{B.9})$$

Applying Lemma A.8, for  $n \gtrsim \log(p-1)$  we have

$$\frac{1}{n} \|(\mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_i)^T \mathbf{X}_{-i}\|_\infty \leq c_7 [\sigma_{x,i,-i} \Lambda_{\max}(\Sigma_{x,-i})]^{1/2} \sqrt{\frac{\log(p-1)}{n}} \quad (\text{B.10})$$

with probability  $\geq 1 - 6c_6 \exp[-(c_7^2 - 1) \log(p-1)]$  for some  $c_6 > 0, c_7 > 1$ . By (B.6), the second term on the right side of (B.9) is bounded above by  $v_\zeta V_x$  with probability  $\geq 1 - 1/p^{\tau_1-2}$  and  $n$  satisfying (B.7). The bound of (A.26) now follows by conditions (C2), (C3) and (B.8).  $\square$

## References

- Belilovsky, E., Varoquaux, G., and Blaschko, M. (2016). Testing for differences in Gaussian graphical models: Applications to brain connectivity. In *NIPS Proceedings*, pages 595–603.
- Bühlmann, P. and van de Geer, S. (2011). *Statistics for High-dimensional Data: Methods, Theory and Applications*. Springer.
- Cai, T. T., Li, H., Liu, W., and Xie, J. (2012). Covariate-adjusted precision matrix estimation with an application in genetical genomics. *Biometrika*, 100(1):139–156.
- Cai, T. T. and Liu, W. (2016). Large-Scale Multiple Testing of Correlations. *J. Amer. Stat. Assoc.*, 111(513):229–240.
- Danaher, P., Wang, P., and Witten, D. M. (2014). The joint graphical lasso for inverse covariance estimation across multiple classes. *J. R. Statist. Soc. B*, 76(2):373–397.
- Gligorijević, V. and Pržulj, N. (2015). Methods for biological data integration: perspectives and challenges. *J. R. Soc. Interface*, 12(112):20150571.
- Gomez-Cabrero, D., Abugessaisa, I., Maier, D., et al. (2014). Data integration in the era of omics: current and future challenges. *BMC Syst. Biol.*, 8(Suppl. 2):11.
- Guo, J., Levina, E., Michailidis, G., and Zhu, J. (2011). Joint estimation of multiple graphical models. *Biometrika*, 98(1):1–15.
- Javanmard, A. and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *J. Mach. Learn. Res.*, 15(1):2869–2909.
- Javanmard, A. and Montanari, A. (2018+). De-biasing the Lasso: Optimal Sample Size for Gaussian Designs. *Ann. Statist.*, To appear. <https://arxiv.org/abs/1508.02757>.
- Joyce, A. R. and Palsson, B. (2006). The model organism as a system: integrating *omics* data sets. *Nat. Rev. Mol. Cell Biol.*, 7:198–210.
- Lee, W. and Liu, Y. (2012). Simultaneous multiple response regression and inverse covariance matrix estimation via penalized Gaussian maximum likelihood. *J. Mult. Anal.*, 111:241–255.

- Li, H., Pouladi, N., Achour, I., et al. (2015). eQTL networks unveil enriched mRNA master integrators downstream of complex disease-associated SNPs. *J. Biomed. Inform.*, 58:226–234.
- Lin, J., Basu, S., Banerjee, M., and Michailidis, G. (2016). Penalized Maximum Likelihood Estimation of Multi-layered Gaussian Graphical Models. *J. Mach. Learn. Res.*, 17:5097–5147.
- Liu, W. (2017). Structural similarity and difference testing on multiple sparse Gaussian graphical models. *Ann. Statist.*, 45(6):2680–2707.
- Liu, W. and Shao, Q.-M. (2014). Phase transition and regularized bootstrap in large-scale  $t$ -tests with false discovery rate control. *Ann. Statist.*, 42(5):2003–2025.
- Ma, J. and Michailidis, G. (2016). Joint Structural Estimation of Multiple Graphical Models. *J. Mach. Learn. Res.*, 17:5777–5824.
- Majumdar, S. and Chatterjee, S. (2018). Non-convex penalized multitask regression using data depth-based penalties. *Stat.*, To appear.
- Mao, Y., Kao, S., Chen, L., et al. (2017). The essential and downstream common proteins of amyotrophic lateral sclerosis: A protein-protein interaction network analysis. *PLoS One*, 12(3):e0172246.
- Meinshausen, N. and Bühlmann, P. (2006). High dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34(3):1436–1462.
- Mitra, R. and Zhang, C.-H. (2016). The benefit of group sparsity in group inference with de-biased scaled group lasso. *Electron. J. Stat.*, 10:1829–1873.
- Obozinski, G., Wainwright, M. J., and Jordan, M. I. (2011). Support Union Recovery in High-dimensional Multivariate Regression. *Ann. Statist.*, 39:1–47.
- Ravikumar, P., Wainwright, M. J., Raskutti, G., and Yu, B. (2011). High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. *Electron. J. Stat.*, 5:935–980.
- Rothman, A. J., Levina, E., and Zhu, J. (2010). Sparse Multivariate Regression With Covariance Estimation. *J. Comp. Graph. Stat.*, 19:947–962.
- Stucky, B. and van de Geer, S. (2017). Asymptotic Confidence Regions for Highdimensional Structured Sparsity. <https://arxiv.org/abs/1706.09231>.
- Tomczak, K., Czerwińska, P., and Wiznerowicz, M. (2015). The Cancer Genome Atlas (TCGA): An immeasurable source of knowledge. *Contemp. Oncol.*, 19:A68–A77.
- van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On Asymptotically Optimal Confidence Regions and Tests for High-Dimensional Models. *Ann. Statist.*, 42:1166–1202.

- Xie, Y., Liu, Y., and Valdar, W. (2016). Joint estimation of multiple dependent Gaussian graphical models with applications to mouse genomics. *Biometrika*, 103(3):493–511.
- Zhang, C.-H. and Zhang, S. S. (2014). Confidence Intervals for Low Dimensional Parameters in High Dimensional Linear Models. *J. R. Statist. Soc. B*, 76:217–242.
- Zhang, Y., Ouyang, Z., and Zhao, H. (2017). A statistical framework for data integration through graphical models with application to cancer genomics. *Ann. Appl. Stat.*, 11(1):161–184.