

Joint Estimation and Inference for Data Integration Problems based on Multiple Multi-layered Gaussian Graphical Models

Subhabrata Majumdar*
and
George Michailidis†
University of Florida, Gainesville, FL, USA

Abstract: The rapid development of high-throughput technologies has enabled the generation of data from biological or disease processes that span multiple layers, like genomic, proteomic or metabolomic data, and further pertain to multiple sources, like disease subtypes or experimental conditions. In this work, we propose a general statistical framework based on Gaussian 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 develop 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 the proposed methodology is evaluated on synthetic data.

Keywords: Data integration; Gaussian Graphical Models; Neighborhood selection; Group lasso; high-dimensional asymptotics; multiple testing; false discovery rate

*Email: smajumdar@ufl.edu

†Corresponding author. Email: gmichail@ufl.edu

1 Introduction

Aberrations in complex biological systems develop in the background of diverse genetic and environmental factors and are associated with multiple complex molecular events. These include changes in the genome, transcriptome, proteome and metabolome, as well as epigenetic effects. Advances in high-throughput profiling techniques have enabled a systematic and comprehensive exploration of the genetic and epigenetic basis of various diseases, including cancer (??), diabetes (??), chronic kidney disease (?), etc. Further, such multi-Omics collections have become available for patients belonging to different, but related disease subtypes, with The Cancer Genome Atlas (TCGA: ?) being a prototypical one. Hence, there is an increasing need for models that can *integrate* such complex data both *vertically* across multiple modalities and *horizontally* across different disease subtypes.

Figure 1 provides a schematic representation of the horizontal and vertical structure of such heterogeneous multi-modal Omics data as outlined above. A simultaneous analysis of all components in this complex layered structure has been coined in the literature 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 (???). A review of the present approaches towards achieving this goal, which are based mostly on specific case studies, can be found in ? and ?.

Contributions are two-fold. Firstly, we propose an integrative framework to conduct simultaneous inference for all parameters in *multiple* and *multi-layer* graphical models, essentially formalizing the structure in Figure 1. We decompose the multi-layer problem into a series of two-layer problems, propose an estimation algorithm for them based on group penalization, and derive theoretical properties of the estimators. Generalizing to group structures on the model parameters allows us to incorporate prior information, *as and when available*, on within-layer or between-layer sub-graph components shared across *some or all* $k = 1, \dots, K$. For biological processes, such information can stem from experimental or mechanistic knowledge (for example a pathway-based grouping of genes). 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. Subsequently, 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.

The novel techniques developed are based on a small number of technical assumptions that are quite general. For example, the model quantities used in our global testing procedure do not necessarily need to be sparse, and instead are only required to have $O(\sqrt{\text{parameter dimension}/n})$ finite sample error bounds that have become standard in the high-dimensional literature (for example see ???). The advantage of this flexibility is that components can be switched out to adapt the framework to other technical assumptions. The optional group sparsity assumptions in our estimation technique can be replaced by other structural restrictions (or no restrictions), for example low-rank or low-rank-plus-

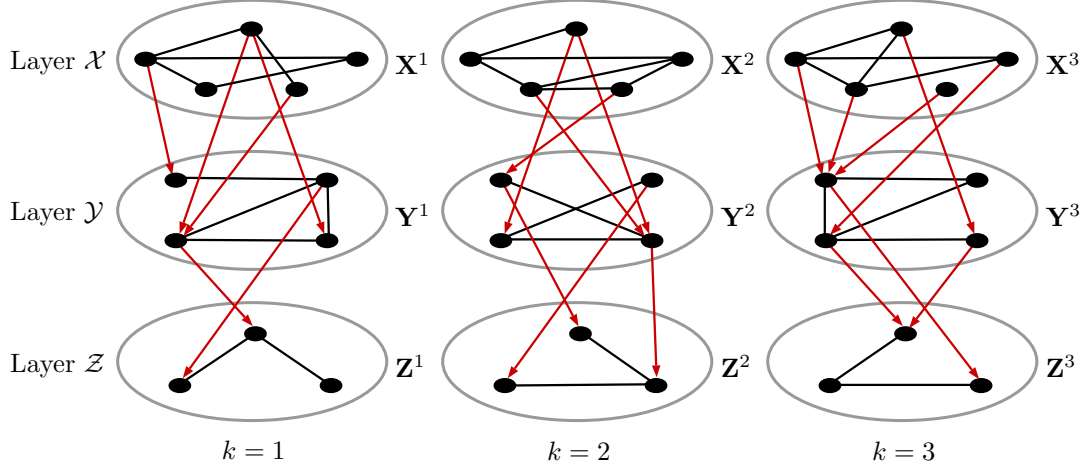


Figure 1: Multiple multilayer graphical models. The matrices $(\mathbf{X}^k, \mathbf{Y}^k, \mathbf{Z}^k)$, $k = 1, 2, 3$ indicate data for each layer and category k . Within-layer connections (black lines) are undirected, while between-layer connections (red lines) go from an upper layer to the successive 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 .

sparse, as deemed appropriate by the prior dependency assumptions across parameters. As long as these resulting estimates converge to the true parameters at the specified finite-sample rates, they can be used by the developed testing methodology.

Related work: 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 from high-dimensional data (e.g. see references in ?), attention has shifted to estimating parameters from more complex structures. This includes (1) analyzing multiple related but not identical graphical models simultaneously, and (2) stacking up more multiple graphical models to form hierarchical multilayer networks, with both directed and undirected edges. For the first class of problems, ? and ? assumed perturbations over a common underlying structure to model multiple precision matrices, while ? proposed using fused/group lasso type penalties for the same task. To incorporate prior information on the group structures across several graphs, ? 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 sets of parameters (??), only recently ? 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.

All the above approaches focus either on the horizontal or the vertical dimensions of

the full hierarchical structure depicted in Figure 1. Hence, multiple related groups of heterogeneous data sets have 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 $(\mathbf{X}^k, \mathbf{Y}^k, \mathbf{Z}^k)$, $k = 1, 2, 3$). In another line of work, ?? model all undirected edges within all nodes together using penalized log-likelihoods. The advantage of this approach is that it can incorporate feedback loops and connections between nodes in non-adjacent layers. However, it has two considerable caveats. Firstly, it does not distinguish between hierarchies, hence delineating the direction of a connection between two nodes across two different Omics modalities is not possible in such models. Secondly, computation becomes difficult when data from different Omics modalities are considered, since the number of estimable parameters increases at a faster rate compared to a hierarchical model.

While there has been some progress for parameter estimation in multilayer models, little is known about the sampling distributions of resulting estimates. Current research on such distributions and related testing procedures for estimates from high-dimensional problems has been limited to single-response regression using lasso (????) or group lasso (?) penalties, and partial correlations of single (?) or multiple (??) GGMs. From a systemic perspective, testing and identifying downstream interactions that differ across experimental conditions or disease subtypes can offer important insights on the underlying biological process (??). In our proposed integrative framework, this can be accomplished by developing a hypothesis testing procedure for entries in the within-layer regression matrices.

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. Sections 4 and ?? are devoted to implementation of our methodology. In Section 4, 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. Section 5 presents a real data example, where we illustrate how the application of our framework leads to knowledge discovery in complex biological networks. We conclude the paper with a discussion in Section ??. Proofs of all theoretical results, as well as some auxiliary results, are given in the Appendix.

Notation: 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 a positive semi-definite matrix \mathbf{P} , we denote its smallest and largest eigenvalues by $\Lambda_{\min}(\mathbf{P})$ and $\Lambda_{\max}(\mathbf{P})$, respectively. 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, ℓ_1 and ℓ_∞ 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 any set \mathcal{S} , $|\mathcal{S}|$ denotes the number of elements in that set. For positive real numbers A, B we write $A \gtrsim B$ 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 data sets, each pertaining to an M -layered Gaussian Graphical Model (GGM) that has p_m nodes in the m -th layer ($1 \leq m \leq M$). The k^{th} model has the following structure:

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

In addition, information is available on horizontal (across k) or vertical (across m) dependencies among nodes within a layer or between nodes of adjacent layers. These are represented by known structured sparsity (i.e. grouping) 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 . Our goal is to leverage this side-information to estimate the full hierarchical structure of the network- specifically to obtain the undirected edges for the nodes inside a single layer, and the directed edges between two successive layers through jointly estimating $\{\Omega_m^k\}$ and $\{\mathbf{B}_m^k\}$.

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

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

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

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

wherein 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 the theoretical discussion in the remainder 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^{th} layer ($m > 1$) and between-layer directed edges from the $(m-1)^{\text{th}}$ layer to the m^{th} layer can be estimated from the corresponding data matrices in a similar fashion (see details in ?). On the other hand, parameters in the very first layer are analogous to $\Omega_x := \{\Omega_x^k\}$, and can be estimated from $\{\mathbf{X}^k\}$ using any method for joint estimation of multiple graphical models (e.g. ??). This provides all building blocks for recovering the full hierarchical structure of our M -layered multiple GGMs.

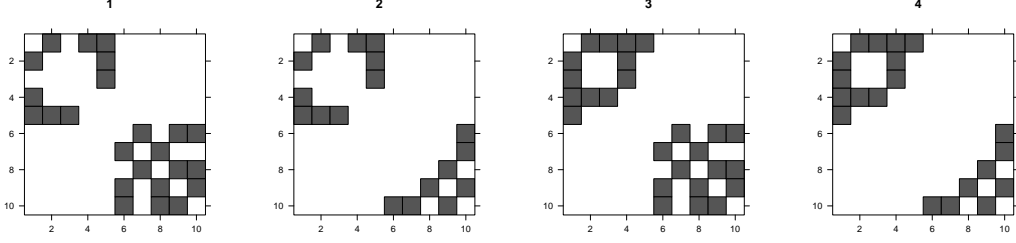


Figure 2: Shared sparsity patterns for four 10×10 precision matrices. for elements $\mathcal{G}_{x,ii'}$ in the upper 5×5 block, matrices (1,2) and (3,4) have the same non-zero support, i.e. $\mathcal{G}_{x,ii'} = \{(1,2), (3,4)\}$. On the other hand, when i, i' are in the lower block, $\mathcal{G}_{x,ii'} = \{(1,3), (2,4)\}$

2.2 Algorithm

We assume an element-wise group sparsity pattern over k for the precision matrices Ω_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 , and consists of non-overlapping index groups g such that $g \subseteq \mathcal{I}_K, \cup_{g \in \mathcal{G}_x^{ii'}} g = \mathcal{I}_K$. First introduced in ?, this formulation helps incorporate group structures that are common across some of the precision matrices being modeled. Figure 2 illustrates this through a small example. Subsequently, we use the Joint Structural Estimation Method (JSEM) (?) to estimate Ω_x , which first uses the group structure given by \mathcal{G}_x in penalized nodewise regressions (?) to obtain neighborhood coefficients $\zeta_i = (\zeta_i^1, \dots, \zeta_i^K)$ of each variable $X_i, i \in \mathcal{I}_p$, then fits a maximum likelihood model over the combined support sets to obtain sparse estimates of the precision matrices:

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

where $\hat{\mathbf{S}}_x^k := (\mathbf{X}^k)^T \mathbf{X}^k / n$, η_n is a tuning parameter, and $\mathbb{S}_+(\hat{E}_x^k)$ is the set of positive-definite matrices that have non-zero supports restricted to \hat{E}_x^k .

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

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

In other words, any arbitrary partition of $\mathcal{I}_p \times \mathcal{I}_q \times \mathcal{I}_K$ can be specified as the sparsity pattern of \mathcal{B} .

Denote the neighborhood coefficients of the j^{th} variable in the lower layer by θ_j^k , and $\Theta_j := (\theta_j^1, \dots, \theta_j^K), \Theta = \{\Theta_j\}$. We obtain sparse estimates of \mathcal{B}, Θ , and subsequently Ω_y ,

by solving the following group-penalized least square minimization problem that has the tuning parameters γ_n and λ_n and then refitting:

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

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

The outcome of a node in the lower layer is thus modeled using all other nodes in that layer using the neighborhood coefficients $\hat{\mathbf{B}}_j^k$, and nodes in the immediate upper layer using the regression coefficients $\hat{\boldsymbol{\theta}}_j^k$.

Remark 1. Common sparsity structures across the same layer are incorporated into the regression by the group penalties over the element-wise groups $\boldsymbol{\theta}_{jj'}^{[g]}$, while sparsity pattern overlaps across the different regression matrices \mathbf{B}^k are handled by the group penalties over $\mathbf{B}^{[h]}$, which denote the collection of elements in \mathcal{B} that are in h . Other kinds of structural assumptions on \mathcal{B} or Θ can be handled within the above structure by swapping out the group norms in favor of other appropriate norm-based penalties.

Remark 2. Group sparsity assumptions are not necessary for the JMMLE framework: rather, they help leverage additional information regarding interaction of features in and between the layers in many applications, *as and when that information is available*. In the vertical direction of the model, i.e. given a fixed k , a framework agnostic of any structural dependency assumptions amounts to element-wise groups in \mathbf{B}^k and Ω_y^k . In JMMLE, this occurs by construction for Θ , and since \mathcal{H} consists of all possible partitions of $\mathcal{I}_p \times \mathcal{I}_q \times \mathcal{I}_K$, it covers the case of element-wise groups as well. On the other hand, the absence of any horizontal (i.e. across k) dependency simply decomposes the problems (2.5) and (2.6) into K independent sub-problems that can be solved separately either by setting $K = 1$ in our framework or by using existing methods, such as ?. The proposed framework provides significant more generality and aims at *tight* vertical and horizontal integration based on available prior information.

2.2.1 Alternating Block Algorithm

The objective function in (2.5) is bi-convex, i.e. convex in \mathcal{B} for fixed Θ , 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 Θ , 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 algorithm as discussed in detail in ?. We choose the initial values $\{\hat{\mathbf{B}}^{k(0)}\}$ by fitting separate lasso

regression models for each j and k :

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

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

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

The steps of our full estimation procedure, coined as the *Joint Multiple Multi-Layer Estimation* (JMMLE) method, are summarized in Algorithm 1.

Algorithm 1. (The JMMLE Algorithm)

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

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

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

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

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

2.2.2 Tuning parameter selection

A number of methods have been proposed in the literature to select regularization tuning parameters in ℓ_1 -penalized problems. Some approaches rely on traditional criteria like cross-validation, Akaike's Information Criterion (AIC) (?) or the Bayesian Information Criterion (BIC) (??). A number of studies have proposed their modifications for the case when feature dimensions increase with sample size (???).

As a demonstration, to select the tuning parameter λ_n we use the High-dimensional BIC (HBIC, ??), and for selecting γ_n in the node-wise regression step in the JSEM model (2.4), employ BIC as in ?. Unlike BIC, the penalty term in HBIC scales with the parameter dimensions. As a result, the tuning parameter selected as the minimizer of HBIC asymptotically identifies the oracle estimator in ultra-high dimensional penalized

problems (??). In our case, we train multiple JMMLE models using Algorithm 1 over a finite set of values $\lambda_n \in \mathcal{D}_n$, and calculate their HBIC:

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

Following this step, we select the optimal λ_n as the empirical minimizer of HBIC over \mathcal{D}_n : $\lambda^* = \arg \min_{\lambda_n \in \mathcal{D}_n} \text{HBIC}(\lambda, \hat{\Theta}_{\gamma_n^*(\lambda_n)})$.

The step for updating Θ , i.e. (2.10), in our JMMLE algorithm is analogous to the JSEM method ?, hence we use BIC to select the penalty parameter γ_n . In our setting the BIC for a given γ_n and fixed \mathcal{B} is given by:

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

where γ_n in subscript indicates the corresponding quantity is calculated taking γ_n 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 γ_n as the one with the smallest BIC over a fixed set of values \mathcal{C}_n . Thus for a fixed $\lambda_n \equiv \lambda$, our final choice of γ_n will be $\gamma_n^*(\lambda) = \arg \min_{\gamma_n \in \mathcal{C}_n} \text{BIC}(\gamma_n; \hat{\mathcal{B}}_{\lambda_n})$.

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 body of the paper, while detailed proofs and auxiliary results are delegated to the Appendix.

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

Definition 1 (Bounded eigenvalues). A positive definite matrix $\Sigma \in \mathbb{M}(b, b)$ is said to have bounded eigenvalues with constants (c_0, d_0) if

$$0 < 1/c_0 \leq \Lambda_{\min}(\Sigma) \leq \Lambda_{\max}(\Sigma) \leq 1/d_0 < \infty$$

Definition 2 (Diagonal dominance). A matrix $\mathbf{M} \in \mathbb{M}(b, b)$ is said to be strictly diagonally dominant if for all $a \in \mathcal{I}_b$,

$$|(\mathbf{M})_{aa}| > \sum_{a' \neq a} |(\mathbf{M})_{aa'}|$$

Also denote $\Delta_0(\mathbf{M}) = \min_a \{ |(\mathbf{M})_{aa}| - \sum_{a' \neq a} |(\mathbf{M})_{aa'}| \}$.

Our first result establishes the convergence of Algorithm 1 for fixed realizations of $(\mathcal{X}, \mathcal{E})$.

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.

As established in Theorems 2.3 and 2.3, at sub-iterations of Algorithm 1 (specifically steps 3 and 4) a $O(\sqrt{\log(pq)}/n)$ bound on $\hat{\mathcal{B}}^{(t)}$ leads to $\hat{\Theta}^{(t+1)}$ being consistent for Θ_0 , and a $O(\sqrt{\log q}/n)$ bound on $\hat{\Theta}^{(t)}$ leads to $\hat{\mathcal{B}}^{(t+1)}$ being consistent for \mathcal{B}_0 : both with probability approaching 1 as $p, q, n \rightarrow \infty$. Thus, while the constant $R(p, q, n)$ in Theorem 2.3 above does not need to obey any explicit bounds for Algorithm 1 to have a limit point, having a tighter $O(\sqrt{\log(pq)}/n)$ bound ensures that the limit point lies close to the population parameters with high enough probability.

The next steps establish that for random realizations of \mathcal{X} and \mathcal{E} ,

- (a) successive iterates lie in this non-expanding ball around the true parameters, and
- (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 sub-problems. Take as $\beta = (\text{vec}(\mathbf{B}^1)^T, \dots, \text{vec}(\mathbf{B}^K)^T)^T$: any subscript or superscript on \mathbf{B} being passed on to β . Denote by $\hat{\Theta}$ and $\hat{\beta}$ the generic estimators given by

$$\hat{\Theta}_j = \arg \min_{\Theta_j \in \mathbb{M}(q-1, K)} \left\{ \frac{1}{n} \sum_{k=1}^K \|\hat{\mathbf{E}}_j^k - \hat{\mathbf{E}}_{-j}^k \theta_j^k\|^2 + \gamma_n \sum_{j \neq j'} \sum_{g \in \mathcal{G}_{jj'}^{jj'}} \|\theta_{jj'}^{[g]}\| \right\}; \quad j \in \mathcal{I}_q, \quad (2.12)$$

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{pqK}} \left\{ -2\beta^T \hat{\gamma} + \beta^T \hat{\Gamma} \beta + \lambda_n \sum_{h \in \mathcal{H}} \|\beta^{[h]}\| \right\}, \quad (2.13)$$

where

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

with

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

Using matrix algebra it is easy to see that solving for β in (2.5) given a fixed $\hat{\Theta}$ is equivalent to solving (2.13).

Next, we assume the following conditions:

- (E1) The matrices $\Omega_{y0}^k, k \in \mathcal{I}_K$ are diagonally dominant,
- (E2) The matrices $\Sigma_{y0}^k, k \in \mathcal{I}_K$ have bounded eigenvalues with constants (c_y, d_y) that are common across k .

Now, we are in a position to establish the estimation consistency for (2.12), as well as the consistency of the final estimates $\hat{\Omega}_y^k$ using their support sets.

Consider random $(\mathcal{X}, \mathcal{E})$, any deterministic $\tilde{\mathcal{B}}$ that satisfy the following bound

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

where C_β depends only on \mathcal{B}_0 . Then, for sample size $n \gtrsim \log(pq)$ there exist constants $c_1, c_3, c_4 > 0, c_2, c_5 > 1, \tau_1 > 2$ such that with probability at least

$$1 - K(1/p^{\tau_1-2} - c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n) - c_4 \exp[-(c_5^2 - 1) \log(pq)]),$$

the following bounds hold:

(I) Denote $|g_{\max}| = \max_{g \in \mathcal{G}_y} |g|$. Then for the choice of tuning parameter

$$\gamma_n \geq 4\sqrt{|g_{\max}|} \mathbb{Q}_0 \sqrt{\frac{\log(pq)}{n}},$$

where \mathbb{Q}_0 depends on the model parameters only, we have

$$\|\hat{\Theta}_j - \Theta_{0,j}\|_F \leq 12\sqrt{s_j} \gamma_n / \psi, \quad (2.15)$$

$$\sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\hat{\boldsymbol{\theta}}_{jj'}^{[g]} - \boldsymbol{\theta}_{0,jj'}^{[g]}\| \leq 48s_j \gamma_n / \psi. \quad (2.16)$$

with $\psi = \min_k \Lambda_{\min}(\Sigma_{x0}^k)/2$.

(II) For the choice of tuning parameter $\gamma_n = 4\sqrt{|g_{\max}|} \mathbb{Q}_0 \sqrt{\log(pq)/n}$,

$$\frac{1}{K} \sum_{k=1}^K \|\hat{\Omega}_y^k - \Omega_{y0}^k\|_F \leq O \left(\mathbb{Q}_0 \sqrt{\frac{|g_{\max}|S}{K}} \sqrt{\frac{\log(pq)}{n}} \right). \quad (2.17)$$

Condition (E2) ensures that the lower layer covariance matrices are well-conditioned, so that the precision matrices $\{\Omega_y^k\}$ exist. The diagonal dominance condition (E1) is a sufficient condition for the convergence bounds of Theorem 2.3 to hold. Specifically, the upper bounds of the finite-sample error rates of estimates $\hat{\Theta}$ and $\hat{\Omega}_y$ are controlled by the ratios of off-diagonal to diagonal elements, i.e. $\omega_{y,jj}^k / \sum_{j' \neq j} |\omega_{y,jj'}^k|$ through the multiplier \mathbb{Q}_0 . A number of ℓ_1 -penalized problems make Restricted Eigenvalue (RE)-type assumptions (???) on the (upper layer) design matrices. Following ?? (see Lemma B.1 and Proposition 1 in respective papers), we utilize the diagonal dominance condition to ensure RE conditions for some key model quantities.

To prove an equivalent result for the solution of (2.13), we need the following conditions on the true parameter versions $(\mathbf{T}_0^k)^2$, defined from Θ_0 similarly as (2.14).

(E3) The matrices $(\mathbf{T}_0^k)^2, k \in \mathcal{I}_K$ are diagonally dominant,

(E4) The matrices $\Sigma_{x0}^k, k \in \mathcal{I}_K$ have bounded eigenvalues with common constants (c_x, d_x) .

Given these, we next establish the required consistency results.

Assume random $(\mathcal{X}, \mathcal{E})$, and fixed $\tilde{\Theta}$ so that for $j \in \mathcal{I}_q$,

$$\|\tilde{\Theta}_j - \Theta_{0,j}\|_F \leq C_\Theta \sqrt{\frac{\log q}{n}}$$

for some C_Θ dependent on Θ_0 only. Then, given the choice of tuning parameter

$$\lambda_n \geq 4\sqrt{|h_{\max}|} \mathbb{R}_0 \sqrt{\frac{\log(pq)}{n}},$$

where \mathbb{R}_0 depends on the population parameters only, with probability at least

$$1 - K(c_1 \exp[-(c_2^2 - 1) \log(pq)] - 2 \exp(-c_3 n))$$

the following bounds hold:

$$\|\hat{\beta} - \beta_0\|_1 \leq 48\sqrt{|h_{\max}|} B \lambda_n / \psi_*, \quad (2.18)$$

$$\|\hat{\beta} - \beta_0\| \leq 12\sqrt{B} \lambda_n / \psi_*, \quad (2.19)$$

$$\sum_{h \in \mathcal{H}} \|\beta^{[h]} - \beta_0^{[h]}\| \leq 48B \lambda_n / \psi_*, \quad (2.20)$$

$$(\hat{\beta} - \beta_0)^T \hat{\Gamma} (\hat{\beta} - \beta_0) \leq 72B \lambda_n^2 / \psi_*, \quad (2.21)$$

where $|h_{\max}| = \max_{h \in \mathcal{H}} |h|$, d_k is the maximum degree $(\mathbf{T}_0^k)^2$, and

$$\psi_* = \frac{1}{2} \min_k \left[\Lambda_{\min}(\Sigma_{x0}^k) \left(\Delta_0((\mathbf{T}_0^k)^2) - d_k C_\Theta \sqrt{\frac{\log(pq)}{n}} \right) \right].$$

Remark 3. In an effort to keep the JMMLE framework as general as possible, we do not impose any explicit sparsity conditions on the fixed quantities $\tilde{\mathcal{B}}$ and $\tilde{\Theta}$ used to estimate the other parameter inside an iteration of Algorithm 1. Since C_β (or C_Θ) depends only on the population parameter \mathcal{B}_0 (or Θ_0), when that parameter is actually sparse their corresponding sparsity values can be a part of C_β (or C_Θ). When we do obtain the actual estimates ((2.15)–(2.17) in Theorem 2.3 and (2.18)–(2.21) in Theorem 2.3), their finite-sample error bounds scale with the corresponding sparsity parameters $\{s_j\}$ and B at rates that are standard in the literature (???)

Following the choice of tuning parameters in Theorems 2.3 and 2.3, $S = o(n/\log(pq))$ and $B = o(n/\log(pq))$ are sufficient conditions on the sparsity of corresponding parameters for the JMMLE estimators to be consistent. As the last step to establish estimation consistency for the limit points of Algorithm 1, we now ensure that the starting values are satisfactory as previously discussed.

Consider the starting values as derived in (2.7) and (2.8). For sample size $n \gtrsim \log(pq)$, and the choice of the tuning parameter

$$\lambda_n \geq 4c_2 \max_{k \in \mathcal{I}_K} \left\{ [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} \right\} \sqrt{\frac{\log(pq)}{n}},$$

we have $\|\hat{\beta}^{(0)} - \beta_0\|_1 \leq 64B\lambda_n/\psi^*$ with probability at least $1 - c_1 \exp(-(c_2^2 - 1)\log(pq)) - 2\exp(-c_3n)$. Also, for $\gamma_n \geq 4\sqrt{|g_{\max}|}\mathbb{Q}_0\sqrt{\log(pq)/n}$ we have

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

with probability at least

$$1 - K(1/p^{\tau_1-2} - c_1 \exp[-(c_2^2 - 1)\log(pq)] - 2\exp(-c_3n) - c_4 \exp[-(c_5^2 - 1)\log(pq)]).$$

Putting all the pieces together, the required consistency result given our choice of starting values follows in a straightforward manner.

Assume conditions (E1)-(E4), 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 λ_n

$$\lambda_n \geq 4 \max \left[c_2 \max_{k \in \mathcal{I}_K} \left\{ [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} \right\}, \sqrt{|h_{\max}|}\mathbb{R}_0 \right] \sqrt{\frac{\log(pq)}{n}},$$

we have

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

with probability at least $1 - 18c_1 \exp[-(c_2^2 - 1)\log(pq)] - 4\exp(-c_3n)$.

(II) For the choice of γ_n

$$\gamma_n \geq 4\sqrt{|g_{\max}|}\mathbb{Q}_0\sqrt{\frac{\log(pq)}{n}},$$

(2.15) and (2.16) hold, while for $\gamma_n = 4\sqrt{|g_{\max}|}\mathbb{Q}_0\sqrt{\log(pq)/n}$, (2.17) holds, both with probability at least

$$1 - K(2/p^{\tau_1-2} - 2c_1 \exp[-(c_2^2 - 1)\log(pq)] - 4\exp(-c_3n) - 2c_4 \exp[-(c_5^2 - 1)\log(pq)]).$$

Remark 4. To save computation time for high data dimensions, an initial screening step, e.g. the debiased lasso procedure of ?, can be used to first restrict the support set of \mathbf{B}_j^k before obtaining the initial estimates using (2.7). The consistency properties of resulting initial and final estimates follow along the lines of the special case $K = 1$ discussed in ?, in conjunction with Theorem 2.3 and Corollary 2.3, respectively. We leave the details to the reader.

Remark 5. While the proof of the above results roughly follow roughly similar roadmaps as the case with $K = 1$ and simple ℓ_1 -penalization (?) and the joint structural estimation of ? and utilize Gaussian concentration inequalities, generalization to $K > 1$ and an optional grouping structure in \mathcal{B} add significant additional technical complexity to the proofs. More importantly, we work in presence of minimal assumptions, steering clear of conditions used in previous works, like Incoherence (?) and Uniform Irrepresentability (?) that are hard to verify in practice.

3 Hypothesis testing in multilayer models

In this section, we lay out a framework for hypothesis testing in our proposed joint multi-layer structure. Present literature in high-dimensional hypothesis testing either focuses on testing for similarities in the within-layer connections of single-layer networks (??), or coefficients of single response penalized regression (???). However, to our knowledge no method is available in the literature to perform testing for *between-layer* connections in a two-layer (or multi-layer) setup.

Denote the i^{th} row of the coefficient matrix \mathbf{B}^k by \mathbf{b}_i^k , for $i \in \mathcal{I}_p$. In this section we are generally interested in obtaining asymptotic sampling distributions of $\widehat{\mathbf{b}}_i^k$, and subsequently formulating testing procedures to detect similarities or differences across k in the full vector \mathbf{b}_i^k or its elements. There are two main challenges in doing the above- firstly the need to mitigate the bias of the group-penalized JMMLE estimators, and secondly the dependency among response nodes translating into the need for controlling false discovery rate while simultaneously testing for several element-wise hypotheses concerning the true values $b_{0ij}^k, j \in \mathcal{I}_q$. To this end, in Section 3.1 we first propose a debiased estimator for \mathbf{b}_i^k that makes use of already computed (using JSEM) node-wise 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, i.e. testing for the null hypothesis $H_0^i : \mathbf{b}_{0i}^1 = \mathbf{b}_{0i}^2$, as well as pairwise simultaneous tests across $j \in \mathcal{I}_q$ for detecting the element-wise differences $b_{0ij}^1 - b_{0ij}^2$.

3.1 Debiased estimators and asymptotic normality

? proposed a debiasing procedure for lasso estimates and subsequently calculate confidence intervals for individual coefficients β_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 $\hat{\beta}^{(\text{init})} \in \mathbb{R}^p$ their debiased estimator was defined as:

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

where \mathbf{z}_j is the vector of residuals from the ℓ_1 -penalized regression of \mathbf{x}_j on \mathbf{X}_{-j} . With centering around the true parameter value, say β_j^0 , and proper scaling this has an asymptotic normal distribution:

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

Essentially, they obtain the debiasing factor for the j^{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. ? later generalized this idea to group lasso estimates. Further, ? and ? 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 $\hat{\mathbf{b}}_i^k$ denotes the i^{th} row of $\hat{\mathbf{B}}^k$, and $t_i^k = (\mathbf{X}_i^k - \mathbf{X}_{-i}^k \hat{\boldsymbol{\zeta}}_i^k)^T \mathbf{X}_i^k / n$, 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 ?. However, as seen 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 through, and they continue to hold for arbitrary sparsity patterns over k in all of 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 ?, ? and ? work on coefficients from lasso regressions, ? debias the coefficients of pre-specified groups in the coefficient vector from a group lasso. Current proposals for obtaining the debiasing factor available in the literature include node-wise lasso (?) and a variance minimization scheme with ℓ_∞ -constraints (?). 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 $k \in \mathcal{I}_K$:

$$\|\hat{\boldsymbol{\zeta}}^k - \boldsymbol{\zeta}_0^k\|_1 \leq D_\zeta = O\left(\sqrt{\frac{\log p}{n}}\right)$$

where D_ζ depends only on the true values, i.e. $\{\boldsymbol{\zeta}_0^k\}$.

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

$$\|\hat{\Omega}_y^k - \Omega_{y0}^k\|_\infty \leq D_\Omega = O\left(\sqrt{\frac{\log(pq)}{n}}\right)$$

where D_Ω depends only on Ω_{y0} .

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

$$\|\widehat{\mathbf{B}}^k - \mathbf{B}_0^k\|_1 \leq D_\beta = O\left(\sqrt{\frac{\log(pq)}{n}}\right),$$

where D_β depends on \mathcal{B}_0 only.

The above finite-sample error rates are common in high-dimensional problems, and can pertain to sparse (???) or non-sparse estimators (??). Based on the estimators plugged in, additional conditions may be involved in the estimation step. For example, JMMLE estimators satisfy these bounds under conditions (E1)-(E4) following the results in Section 2.3.

Given these conditions, the following result 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 ?. However, the authors use the nuclear norm as the loss function while obtaining the debiasing factors and employ 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 random design matrices.

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

$$\begin{aligned}\widehat{\Omega}_y &= \text{diag}(\widehat{\Omega}_y^1, \dots, \widehat{\Omega}_y^K) \\ \mathbf{M}_i &= \text{diag}(m_i^1, \dots, m_i^K) \\ \widehat{\mathbf{C}}_i &= \text{vec}(\widehat{\mathbf{c}}_i^1, \dots, \widehat{\mathbf{c}}_i^K)^T \\ \mathbf{D}_i &= \text{vec}(\mathbf{b}_{0i}^1, \dots, \mathbf{b}_{0i}^K)^T\end{aligned}$$

Also assume that conditions (E2), (E4) hold, and the matrices $\Omega_{x0}^k, k \in \mathcal{I}_K$ are diagonally dominant. Then, for sample size satisfying $\log p = o(n^{1/2})$, $\log q = o(n^{1/2})$ we have

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

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

3.2 Test formulation

We now simply plug in estimators from the JMMLE algorithm in Theorem 3.1. Doing so is fairly straightforward. Condition (T1) is ensured by the JSEM penalized neighborhood estimators in (2.4) (immediate from Proposition A.1 in ?). On the other hand, bounds on total sparsity of the true coefficient matrices: $B = o(\sqrt{n}/\log(pq))$, and lower layer precision matrices: $S = o(n/\log(pq))$, in conjunction with Corollary 2.3, ensure conditions (T2) and (T3), respectively -all with probability approaching 1 as $(n, p, q) \rightarrow \infty$.

An asymptotic joint distribution of debiased versions of the JMMLE regression estimates can then be obtained immediately. Consider the estimates $\widehat{\mathbf{B}}$ and $\widehat{\Omega}_y$ obtained from Algorithm 1, and upper layer neighborhood coefficients from solving the node-wise regression in (2.4). Suppose that $\log(pq)/\sqrt{n} \rightarrow 0$, and the sparsity conditions $B = o(\sqrt{n}/\log(pq))$, $S = o(n/\log(pq))$ are satisfied. Then, with the same notations as in Theorem 3.1 we have

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

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 individual rows of the true coefficient matrices, i.e. \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 multivariate normal distribution. Consequently, we formulate a global test for detecting differential overall downstream effect of the i^{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, the above testing procedure maintains rate optimal power.

Consider the global test given in Algorithm 2, performed using parameter estimates satisfying conditions (T1)-(T3). Define $\boldsymbol{\delta} := \mathbf{b}_{0i}^1 - \mathbf{b}_{0i}^2$. Further, assume that either of the following sufficient conditions are satisfied.

- (I) The following bound holds: $D_\Omega \leq \Delta_0(\Omega_{y0}^k), k \in \mathcal{I}_K$;
- (II) For every $j \in \mathcal{I}_q, k \in \mathcal{I}_K$, we have $\sum_{j'=1}^q |\sigma_{y0, jj'}^k|^q \leq c_0(p)$ for some $q \in [0, 1)$ and positive-valued function $c_0(\cdot)$.

Denote $\sigma_{x0, i, -i}^k = \text{Var}(X_i^k - \mathbb{X}_{-i}^k \boldsymbol{\zeta}_{0, i}^k)$. Then, the power of the global test is given by

$$K_q \left(\chi_{q, 1-\alpha}^2 + n \boldsymbol{\delta}^T \left(\frac{\Sigma_{y0}^1}{\sigma_{x0, i, -i}^1} + \frac{\Sigma_{y0}^2}{\sigma_{x0, i, -i}^2} \right)^{-1} \boldsymbol{\delta} \right) + o(1)$$

where K_q is the cumulative distribution function of the χ_q^2 distribution. Consequently, for $\|\delta\| > O(n^{-1/2})$, $P(H_0^i \text{ is rejected}) \rightarrow 1$ as $(n, p, q) \rightarrow \infty$.

The conditions (I) or (II) above are needed to derive upper bounds for $\|\hat{\Sigma}_y^k - \Sigma_{y0}^k\|_\infty$ using those for $\|\hat{\Omega}_y^k - \Omega_{y0}^k\|_\infty$. While (I) imposes a potentially more stringent bound on the estimation error of Ω_y , (II) restricts the power calculations to a uniformity class of covariance matrices (??).

Remark 6. While the formulation of the testing procedure broadly gives parallel results as ? and ?, it does so *without assuming any specific penalty function* or (group) sparsity conditions (such as strong group sparsity in ?). Instead, we only require the standard finite-sample bounds (T1)-(T3) that are satisfied by both sparse and non-sparse estimators in a high-dimensional setting.

3.3 Control of False Discovery Rate

Given that 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^{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 τ , 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 α , 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 $|d_{ij}| \geq \hat{\tau}$.

To ensure that this procedure maintains 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 [?](#), we consider the following two types of dependencies:

(D1) Define $r_{jj'}^k = \sigma_{y0,jj'}^k / \sqrt{\sigma_{y0,jj}^k \sigma_{y0,j'j'}^k}$ for $j, j' \in \mathcal{I}_q, k = 1, 2$. Suppose there exists $0 < r < 1$ such that $\max_{1 \leq j < j' \leq q} |r_{jj'}^k| \leq r$, and for every $j \in \mathcal{I}_q$,

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

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

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

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

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

Originally proposed by [?](#), the above dependency conditions are meant to control the amount of correlation amongst 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 $q \rightarrow \infty$.

Suppose $\mu_j = b_{0,ij}^1 - b_{0,ij}^2, \sigma_j^2 = \sigma_{y0,jj}^1 / \sigma_{x0,i,-i}^1 + \sigma_{y0,jj}^2 / \sigma_{x0,i,-i}^2$. Assume the following holds as $(n, q) \rightarrow \infty$

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

Next, consider 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 ([??](#)).

Remark 7. Based on the FDR control procedure in Algorithm [3](#), we can perform *within-row thresholding* in the matrices $\hat{\mathbf{B}}^k$ to tackle group misspecification.

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

Even without group misspecification, this helps identify directed edges between layers that have high nonzero values. Similar post-estimation thresholdings have been proposed in the context of multitask regression (??) and neighborhood selection (?). However, our procedure is the first one to provide explicit guarantees on the amount of false discoveries while doing so.

Remark 8. Following (3.5), a sufficient condition on the sparsity of \mathcal{B}_0 for FDR to be asymptotically controlled at some specified level is $B = o(n^\zeta / \log q)$ if (D1) is satisfied, and $B = o(n^{1/3} / \log q)$ if (D1*) is satisfied. In comparison, our results for the global testing procedure require $B = o(\sqrt{n} / \log(pq))$, and point estimation requires $B = o(n / \log(pq))$. In finite samples settings, the stricter sparsity requirements translate to higher sample sizes being needed (given the same (p, q)) for our testing procedures to have satisfactory performances compared to estimation only (See Sections 4.1 and 4.2).

In recent work, ? showed that the $o(\sqrt{n} / \log p)$ bound on the sparsity of the true coefficient vector required to construct confidence intervals from debiased lasso coefficient estimates (???) can be weakened to $o(n / (\log p)^2)$ when the random design precision matrix is known, or is unknown but satisfies certain sparsity assumptions. Similar relaxations may be possible in our case. For example, the machinery in ?, which performs simultaneous testing in multiple (single layer) GGMs using slightly modified FDR thresholds, can be useful in obtaining (3.6) for $\log q = o(n^{1/2})$ under (D1), (D1*) or other suitable dependency assumptions.

3.4 Effect of tuning parameter selection

A topic not adequately addressed in the high-dimensional hypothesis testing literature concerns the effect of the regularization tuning parameter selection methods (HBIC for λ_n and BIC for γ_n in our case) on the size, power and confidence intervals obtained. Ideally, tuning parameter selection method(s) in the estimation step should ensure that the estimated quantities from the model with optimal tuning parameter choices can be plugged into the debiasing procedure to obtain quantities that obey the correct asymptotic properties and are used in the tests that follow (e.g. Algorithms 2 and 3).

In our case, the broad-based assumptions (T1)-(T3) allow plugging in estimators with finite-sample error rates that are satisfied by a host of high-dimensional methods, as previously discussed. Given that the tuning parameters λ_n and γ_n are selected to be above thresholds that scale with feature and sample dimensions, JMMLE estimators adhere to these error rates with high probability (Corollary 2.3), ensuring the correctness of our testing procedures. To empirically make it likely that the correct tuning parameters get selected, in our numerical examples (on synthetic and real data) that follow, we obtain JMMLE estimates over ranges of λ_n and γ_n that scale with the error rates of estimators. Additional technicalities will be involved for a more rigorous analysis, possibly using technical material from approaches such as ??. We defer this topic to future work.

4 Numerical experiments

Next, 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, respectively), and also introduce some computational techniques that significantly accelerate calculations 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) integration 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 ?. 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 gives 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 3 illustrates this structure. We set an off-diagonal element inside each of these common blocks (i.e. A, B, C and α, β 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 Ω_{x0}^k are generated by putting together the corresponding common blocks, their positive definiteness ensured by setting all diagonal elements to be $1 + |\Lambda_{\min}(\Omega_{x0}^k)|$. Then, we get elements in the covariance matrix as

$$\sigma_{x0,ii'}^k = (\bar{\Omega}_{x0}^k)_{ii'} / \sqrt{(\bar{\Omega}_{x0}^k)_{ii}(\bar{\Omega}_{x0}^k)_{i'i'}}, \text{ where } \bar{\Omega}_{x0}^k = (\Omega_{x0}^k)^{-1},$$

and generate rows of \mathbf{X}^k independently from $\mathcal{N}(0, \Sigma_{x0}^k)$. We obtain Σ_{y0}^k and then \mathbf{E}^k using the same setup but with the number of variables being q and setting off-diagonal 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 $b_{0,ij}^k$ non-zero across all k with probability $\pi \in \{5/p, 30/p\}$, 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 150 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 scaled arrays of tuning parameters to train Algorithm 1-

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

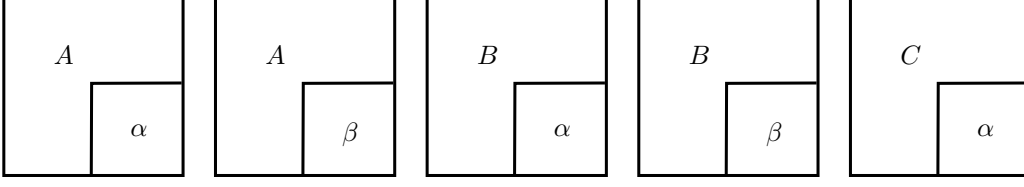


Figure 3: Shared sparsity patterns across k for the precision matrices $\{\Omega_{x0}^k\}$ and $\{\Omega_{y0}^k\}$

using a one-step version of the algorithm (Section 4.3) to save computation time.

We use the following performance metrics to evaluate our estimates $\hat{\mathcal{B}} = \{\hat{\mathbf{B}}^k\}$:

- True positive Rate-

$$\text{TPR}(\hat{\mathbf{B}}_k) = \frac{|\text{supp}(\hat{\mathbf{B}}^k) \cap \text{supp}(\mathbf{B}_0^k)|}{|\text{supp}(\mathbf{B}_0^k)|}; \quad \text{TPR}(\hat{\mathcal{B}}) = \frac{1}{K} \sum_{k=1}^K \text{TPR}(\hat{\mathbf{B}}_k).$$

- True negative Rate-

$$\text{TNR}(\hat{\mathbf{B}}_k) = \frac{|\text{supp}^c(\hat{\mathbf{B}}^k) \cap \text{supp}^c(\mathbf{B}_0^k)|}{|\text{supp}^c(\mathbf{B}_0^k)|}; \quad \text{TNR}(\hat{\mathcal{B}}) = \frac{1}{K} \sum_{k=1}^K \text{TNR}(\hat{\mathbf{B}}_k).$$

- Matthews Correlation Coefficient-

$$\text{TP}(\hat{\mathbf{B}}_k) = |\text{supp}(\hat{\mathbf{B}}^k) \cap \text{supp}(\mathbf{B}_0^k)|; \quad \text{TN}(\hat{\mathbf{B}}_k) = |\text{supp}^c(\hat{\mathbf{B}}^k) \cap \text{supp}^c(\mathbf{B}_0^k)|,$$

$$\text{FP}(\hat{\mathbf{B}}_k) = |\text{supp}^c(\mathbf{B}_0^k)| - \text{TN}(\hat{\mathbf{B}}_k); \quad \text{FN}(\hat{\mathbf{B}}_k) = |\text{supp}(\mathbf{B}_0^k)| - \text{TP}(\hat{\mathbf{B}}_k),$$

$$\text{MCC}(\hat{\mathbf{B}}_k) =$$

$$\frac{\text{TP}(\hat{\mathbf{B}}_k)\text{TN}(\hat{\mathbf{B}}_k) - \text{FP}(\hat{\mathbf{B}}_k)\text{FN}(\hat{\mathbf{B}}_k)}{\sqrt{(\text{TP}(\hat{\mathbf{B}}_k) + \text{FP}(\hat{\mathbf{B}}_k))(\text{TP}(\hat{\mathbf{B}}_k) + \text{FN}(\hat{\mathbf{B}}_k))(\text{TN}(\hat{\mathbf{B}}_k) + \text{FP}(\hat{\mathbf{B}}_k))(\text{TN}(\hat{\mathbf{B}}_k) + \text{FN}(\hat{\mathbf{B}}_k))}},$$

$$\text{MCC}(\hat{\mathcal{B}}) = \frac{1}{K} \sum_{k=1}^K \text{MCC}(\hat{\mathbf{B}}_k).$$

- Relative error in Frobenius norm-

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

We use the same metrics to evaluate the precision matrix estimates $\hat{\Omega}_y^k$ as well, with TPR and TNR calculations confined to off-diagonal entries.

Tables 1 and 2 summarize the results. For estimation of \mathcal{B} , we compare our results to the method in ? that estimates parameters in each of the K two-layer structure separately,

(π_x, π_y)	(p, q, n)	Method	TPR	TNR	MCC	RF
$(5/p, 5/q)$	$(60, 30, 100)$	JMMLE	0.97(0.02)	0.99(0.003)	0.96(0.014)	0.24(0.033)
		Separate	0.96(0.018)	0.99(0.004)	0.93(0.014)	0.22(0.029)
	$(30, 60, 100)$	JMMLE	0.97(0.013)	0.99(0.002)	0.96(0.008)	0.27(0.024)
		Separate	0.99(0.009)	0.99(0.003)	0.93(0.017)	0.18(0.021)
	$(200, 200, 150)$	JMMLE	0.98(0.011)	1.0(0)	0.99(0.005)	0.16(0.025)
		Separate	0.99(0.001)	0.99(0.001)	0.88(0.009)	0.18(0.007)
	$(300, 300, 150)$	JMMLE	1.0(0.001)	1.0(0)	0.99(0.001)	0.14(0.015)
		Separate	1.0(0.001)	0.99(0.001)	0.84(0.01)	0.21(0.007)
	$(200, 200, 100)$	JMMLE	0.97(0.017)	1.0(0)	0.98(0.008)	0.21(0.032)
		Separate	0.32(0.01)	0.99(0.001)	0.49(0.009)	0.85(0.06)
$(30/p, 30/q)$	$(200, 200, 200)$	JMMLE	0.99(0.006)	1.0(0)	0.99(0.007)	0.13(0.016)
		Separate	0.97(0.004)	0.98(0.001)	0.93(0.002)	0.19(0.07)

Table 1: Table of outputs for estimation of regression matrices, giving empirical mean and standard deviation (in brackets) of each evaluation metric over 150 replications.

(π_x, π_y)	(p, q, n)	Method	TPR	TNR	MCC	RF
$(5/p, 5/q)$	$(60, 30, 100)$	JMMLE	0.76(0.018)	0.90(0.006)	0.61(0.024)	0.32(0.008)
		Separate	0.77(0.031)	0.92(0.007)	0.56(0.03)	0.51(0.017)
		JSEM	0.24(0.013)	0.8(0.003)	0.05(0.015)	1.03(0.002)
	$(30, 60, 100)$	JMMLE	0.7(0.018)	0.94(0.002)	0.55(0.018)	0.3(0.005)
		Separate	0.76(0.041)	0.89(0.015)	0.59(0.039)	0.49(0.014)
		JSEM	0.13(0.005)	0.9(0.001)	0.03(0.007)	1.04(0.001)
	$(200, 200, 150)$	JMMLE	0.68(0.017)	0.98(0)	0.48(0.013)	0.26(0.002)
		Separate	0.78(0.019)	0.97(0.001)	0.55(0.012)	0.6(0.007)
		JSEM	0.05(0.002)	0.97(0)	0.02(0.002)	1.01(0)
	$(300, 300, 150)$	JMMLE	0.71(0.014)	0.98(0)	0.44(0.008)	0.25(0.002)
		Separate	0.71(0.017)	0.98(0.001)	0.51(0.011)	0.59(0.005)
		JSEM	0.04(0.002)	0.98(0)	0.02(0.002)	1.01(0)
	$(200, 200, 100)$	JMMLE	0.77(0.016)	0.98(0)	0.46(0.013)	0.31(0.003)
		Separate	0.57(0.027)	0.44(0.007)	0.04(0.008)	0.84(0.002)
		JSEM	0.05(0.002)	0.97(0)	0.01(0.002)	1.01(0)
$(30/p, 30/q)$	$(200, 200, 200)$	JMMLE	0.76(0.018)	0.98(0)	0.55(0.015)	0.27(0.004)
		Separate	0.73(0.023)	0.94(0.003)	0.39(0.017)	0.62(0.011)
		JSEM	0.05(0.002)	0.97(0)	0.03(0.003)	1.01(0)

Table 2: Table of outputs for estimation of lower layer precision matrices over 150 replications.

while for estimation of Ω_y , we compare them with the results in ? and using the single-layer JSEM (?) that estimates Ω_y assuming structured sparsity patterns and centered matrices \mathbf{Y}^k , but not the data in the upper layer, i.e. \mathcal{X} .

Our joint method has higher average MCC across all data settings than the separate method for the estimation of \mathcal{B} , although TPR and TNR values are similar, except for $p = 200, q = 200, n = 100$ where JMMLE has a much higher average TPR. For estimation

(π_x, π_y)	(p, q, n)	TPR($\hat{\mathcal{B}}$)	TNR($\hat{\mathcal{B}}$)	MCC($\hat{\mathcal{B}}$)	RF($\hat{\mathcal{B}}$)
$(5/p, 5/q)$	(60,30,100)	0.98 (0.01)	0.99 (0.002)	0.89 (0.017)	0.29 (0.014)
	(30,60,100)	0.94 (0.022)	0.99 (0.003)	0.93 (0.016)	0.31 (0.028)
	(200,200,150)	0.99 (0.002)	0.99 (0)	0.98 (0.004)	0.17 (0.007)
	(300,300,150)	0.99 (0.001)	1 (0)	0.99 (0.002)	0.15 (0.006)
$(30/p, 30/q)$	(200,200,100)	0.99 (0.006)	1 (0)	0.98 (0.005)	0.2 (0.014)
	(200,200,200)	0.99 (0.009)	1 (0)	0.98 (0.005)	0.15 (0.017)
(π_x, π_y)	(p, q, n)	TPR($\hat{\Omega}_y$)	TNR($\hat{\Omega}_y$)	MCC($\hat{\Omega}_y$)	RF($\hat{\Omega}_y$)
$(5/p, 5/q)$	(60,30,100)	0.71 (0.024)	0.90 (0.005)	0.64 (0.024)	0.34 (0.008)
	(30,60,100)	0.7 (0.019)	0.94 (0.002)	0.59 (0.014)	0.3 (0.004)
	(200,200,150)	0.62 (0.012)	0.98 (0)	0.43 (0.009)	0.27 (0.003)
	(300,300,150)	0.69 (0.013)	0.98 (0)	0.39 (0.008)	0.26 (0.02)
$(30/p, 30/q)$	(200,200,100)	0.78 (0.024)	0.98 (0)	0.43 (0.012)	0.31 (0.003)
	(200,200,200)	0.69 (0.026)	0.98 (0.001)	0.5 (0.02)	0.29 (0.004)

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

(π_x, π_y)	(p, q, n)	FDR
$(5/p, 5/q)$	(60,30,100)	0.19 (0.077)
	(30,60,100)	0.08 (0.064)
	(200,200,150)	0.04 (0.016)
	(300,300,150)	0.02 (0.007)
$(30/p, 30/q)$	(200,200,100)	0.03 (0.019)
	(200,200,200)	0.03 (0.016)

Table 4: Table of outputs giving empirical FDR for estimating \mathcal{B} using JMMLE in presence of group misspecification

of Ω_y , incorporating information from the upper layer vastly improves performance, as demonstrated by the differences in performance between JMMLE and JSEM. For the 4 data settings with lower sparsity ($\pi_x = \pi = 5/p, \pi_y = 5/q$), JMMLE produces sparser estimates compared to the separate method while estimating Ω_y as is evident from the lower TPR and MCC values. However, the RF values indicate that the quality of JMMLE estimates is significantly better. This in fact is a common pattern across the estimation of both \mathcal{B} and Ω_y : JMMLE gives more accurate estimates across the methods, with lower average RF values across all data settings. Finally, for the estimation of Ω_y , JMMLE does better in both of the higher sparsity settings ($\pi_x = \pi = 30/p, \pi_y = 30/q$) across all metrics.

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 task, we first set individual elements inside a non-zero group to be zero with probability 0.2 while generating the data, then pass the JMMLE estimates $\hat{\mathbf{B}}^k$ through the FDR controlling thresholds as given in (3.7). The results are

(π_x, π_y)	(p, q, n)	Method	Global test		Simultaneous test	
			Power	Size	Power	FDR
$(5/p, 5/q)$	(60,30,100)	JMMLE	0.98 (0.016)	0.07 (0.011)	0.94 (0.023)	0.24 (0.027)
		Separate	0.99 (0.007)	0.12 (0.02)	0.91 (0.025)	0.34(0.038)
		SepLasso	0.99 (0.007)	0.11 (0.02)	0.91 (0.025)	0.33(0.038)
	(60,30,200)	JMMLE	0.99 (0.014)	0.07 (0.014)	0.97 (0.013)	0.22 (0.032)
		Separate	0.99 (0.005)	0.08 (0.014)	0.94 (0.019)	0.26(0.031)
		SepLasso	0.99 (0.004)	0.08 (0.014)	0.94 (0.019)	0.26(0.033)
	(30,60,100)	JMMLE	0.98 (0.024)	0.07 (0.014)	0.92 (0.027)	0.24 (0.035)
		Separate	1 (0)	0.07 (0.015)	0.86 (0.036)	0.25(0.039)
		SepLasso	1 (0)	0.08 (0.014)	0.85 (0.036)	0.25(0.039)
	(30,60,200)	JMMLE	0.99 (0.019)	0.08 (0.016)	0.96 (0.023)	0.24 (0.038)
		Separate	1 (0)	0.06 (0.013)	0.9 (0.038)	0.21(0.035)
		SepLasso	1 (0)	0.06 (0.012)	0.91 (0.038)	0.21(0.034)
	(200,200,150)	JMMLE	0.99 (0.006)	0.06 (0.003)	0.84 (0.011)	0.22 (0.007)
		Separate	1 (0)	0.2 (0.008)	0.93 (0.006)	0.46(0.009)
		SepLasso	1 (0)	0.2 (0.008)	0.93 (0.006)	0.46(0.009)
	(300,300,150)	JMMLE	0.99 (0.004)	0.07 (0.009)	0.54 (0.031)	0.34 (0.016)
		Separate	1 (0)	0.27 (0.01)	0.79 (0.007)	0.58(0.008)
		SepLasso	1 (0)	0.27 (0.01)	0.79 (0.007)	0.58(0.008)
	(300,300,300)	JMMLE	0.99 (0.003)	0.03 (0.002)	0.99 (0.003)	0.12 (0.006)
		Separate	1 (0)	0.16 (0.005)	0.99 (0.004)	0.4 (0.007)
		SepLasso	1 (0)	0.16 (0.005)	0.99 (0.004)	0.4 (0.007)
$(30/p, 30/q)$	(200,200,100)	JMMLE	0.99 (0.005)	0.112 (0.003)	0.41 (0.008)	0.52 (0.007)
		Separate	1 (0)	0.47 (0.008)	0.75 (0.007)	0.71(0.004)
		SepLasso	1 (0)	0.47 (0.008)	0.75 (0.007)	0.71(0.004)
	(200,200,200)	JMMLE	0.99 (0.004)	0.09 (0.004)	0.96 (0.006)	0.27 (0.008)
		Separate	1 (0)	0.42 (0.011)	0.98 (0.005)	0.63(0.006)
		SepLasso	1 (0)	0.42 (0.011)	0.98 (0.005)	0.63(0.006)
	(200,200,300)	JMMLE	0.99 (0.002)	0.06 (0.003)	0.99 (0.004)	0.19 (0.008)
		Separate	1 (0)	0.27 (0.01)	0.99 (0.004)	0.52 (0.009)
		SepLasso	1 (0)	0.27 (0.01)	0.99 (0.004)	0.52 (0.009)

Table 5: Table of outputs for global and simultaneous hypothesis testing.

summarized in Tables 3 and 4. Across the simulation settings, values of all metrics are very close to the correctly specified counterparts in Table 1. Thus, the thresholding step proves largely effective. Also, in all cases the empirical FDR for estimating entries in \mathcal{B} is below 0.2. The performance is slightly worse than the correctly specified cases when estimating Ω_y . This is expected, as the estimates $\hat{\Omega}_y$ are obtained from neighborhood coefficients that are calculated based on the *pre-thresholding* coefficient estimates.

4.2 Simulation 2: testing

We slightly change the 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 π , 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} , where $(\mathbf{D})_{ij}, i \in \mathcal{I}_p, j \in \mathcal{I}_q$ takes 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 identical sparsity structures for the pairs of precision matrices $\{\Omega_{x0}^1, \Omega_{x0}^2\}$ and $\{\Omega_{y0}^1, \Omega_{y0}^2\}$. We use 150 replications of the above setup to calculate empirical power of global tests, as well as empirical power and FDR of simultaneous tests. To get the empirical sizes of global tests we use estimators obtained from applying JMMLE on a separate set of data generated setting all elements of \mathbf{D} to 0. The type-I error of global tests is controlled at level 0.05, while FDR is set at 0.2 obtained by calculating the respective thresholds.

Table 5 reports the empirical mean and standard deviations (in brackets) of all relevant quantities computed from debiased coefficients obtained from JMMLE, separate estimation, as well as from applying the original debiasing technique of ? on qK separate lasso estimates of row-level coefficient vectors, i.e. $\hat{\mathbf{b}}_i^k$. 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 from the theoretical analysis, higher sample sizes than those used in Section 4.1 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. While separate estimation has slightly higher power in global testing, our joint method gives better results everywhere else. The empirical size of the JMMLE-based global tests remain slightly higher than the nominal level across the settings considered. This is in all likelihood a consequence of the higher sample size requirements in testing than estimation, as nominal sizes for JMMLE estimates tend to go down when p, q are kept constant and n is increased (for example, check setting 6 vs. setting 7 ($p = 300, q = 300, n = 150, 300$), and setting 8 vs. 9 vs. 10 ($p = 200, q = 200, n = 100, 200, 300$)). This pattern is similar to empirical results in past proposals of high-dimensional testing methodology (??). However, the size estimates of JMMLE are much closer to the nominal level across data settings compared to either separate estimation or separate lasso, owing the fact that only JMMLE is able to leverage information across the different multilayer networks.

4.3 Computation

Next, we discuss some observations and strategies that speed up the JMMLE algorithm and reduce computation 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 $K = 1$ (?), 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:

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

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

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

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 (< 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 *inside each iteration* after the matrices $\widehat{\mathbf{B}}^k$ are updated:

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

$$\widehat{\mathbf{B}}_j^{k(t+1)} = \left[(\mathbf{X}_{\mathcal{S}_{jk}}^k)^T (\mathbf{X}_{\mathcal{S}_{jk}}^k) \right]^{-1} (\mathbf{X}_{\mathcal{S}_{jk}}^k)^T \mathbf{Y}_j^k$$

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

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 $\{\widehat{\Theta}_j\}$ in each iteration. In practice, the algorithm speeds up significantly without compromising on estimation accuracy if we dispense of the Θ update 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 $\widehat{\mathcal{B}}$ using (2.7).
2. Initialize $\widehat{\Theta}$ using (2.8).
3. Update $\widehat{\mathcal{B}}$ as:

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

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

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

(p, q, n)	Method	TPR($\hat{\mathcal{B}}$)	TNR($\hat{\mathcal{B}}$)	MCC($\hat{\mathcal{B}}$)	RF($\hat{\mathcal{B}}$)
(60,30,100)	Full	0.982 (0.013)	0.994 (0.003)	0.959 (0.014)	0.23 (0.021)
	One step	0.971 (0.02)	0.996 (0.003)	0.965 (0.014)	0.242 (0.033)
(30,60,100)	Full	0.966 (0.015)	0.991 (0.003)	0.954 (0.008)	0.269 (0.026)
	One step	0.968 (0.013)	0.992 (0.002)	0.957 (0.008)	0.265 (0.024)
(p, q, n)	Method	TPR($\hat{\Omega}_y$)	TNR($\hat{\Omega}_y$)	MCC($\hat{\Omega}_y$)	RF($\hat{\Omega}_y$)
(60,30,100)	Full	0.756 (0.019)	0.907 (0.005)	0.616 (0.021)	0.318 (0.007)
	One step	0.764 (0.018)	0.904 (0.006)	0.678 (0.024)	0.321 (0.008)
(30,60,100)	Full	0.695 (0.016)	0.943 (0.002)	0.552 (0.015)	0.304 (0.005)
	One step	0.696 (0.018)	0.943 (0.002)	0.552 (0.018)	0.304 (0.005)

Table 6: Comparison of evaluation metrics for full and one-step versions of the JMMLE algorithm.

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

Compared to one-step algorithms based on first order approximation of the objective function (??), we let \mathcal{B} converge completely, then use these solutions to recover the support set of the precision matrices. The estimation accuracy of Ω_y depends on the solution $\hat{\mathcal{B}}$ used to solve the sub-problem (2.12) (Theorem 2.3 and Lemmas A and A). Thus, letting \mathcal{B} converge first ensures that the solutions $\hat{\Theta}$ and $\hat{\Omega}_y$ obtained subsequently are of a better quality compared to a simple early stopping of the JMMLE algorithm.

(p, q, n)	Method	Comp. time (min)
(60,30,100)	Full	6.1
	One-step	0.7
(30,60,100)	Full	22.4
	One-step	2.7

Table 7: Comparison of computation times (averaged over 150 replications) for full and one-step versions of the JMMLE algorithm.

We compared the performance of both versions of our algorithm for the two data settings with smaller feature dimensions. Computations were performed on the HiperGator supercomputer¹, in parallel across 8 cores of an Intel E5-2698v3 2.3GHz processor with 2GB RAM per core, the parallelization being done across the range of values for λ_n within each replication. As seen in Table 6, performance is indistinguishable across all the metrics, but the one-step algorithm saves a significant amount of computation time compared to the full version (Table 7).

5 Discussion

This work introduces an integrative framework for knowledge discovery in multiple multi-layer Gaussian Graphical Models. We exploit *a priori* known structural similarities across

¹<https://www.rc.ufl.edu/services/hipergator>

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.

5.1 Performance improvement

The JMMLE algorithm due to the incorporation of prior information about sparsity patterns improves on the theoretical convergence rates of the estimation method for *single* multi-layer GGMs (i.e. $K = 1$) introduced in ?. With our initial estimates, the method of ? achieves the following convergence rates for the estimation of \mathcal{B} and Ω_y , respectively (using Corollary 4 therein):

$$\begin{aligned}\|\hat{\beta} - \beta_0\|_F &\leq \sum_{k=1}^K O\left(\sqrt{\frac{b_k \log(pq)}{n}}\right), \\ \sum_{k=1}^K \|\hat{\Omega}_y^k - \Omega_{y0}^k\|_F &\leq O\left(K\sqrt{\frac{(S+q)\log(pq)}{n}}\right).\end{aligned}$$

In comparison, JMMLE has the following rates:

$$\begin{aligned}\|\hat{\beta} - \beta_0\|_F &\leq O\left(\sqrt{\frac{|h_{\max}|B\log(pq)}{n}}\right), \\ \sum_{k=1}^K \|\hat{\Omega}_y^k - \Omega_{y0}^k\|_F &\leq O\left(\sqrt{\frac{KS|g_{\max}|\log(pq)}{n}}\right).\end{aligned}$$

For \mathcal{B} , joint estimation outperforms separate estimation when group sizes are small, so that $(|h_{\max}|B)^{1/2} < \sum_k b_k^{1/2}$. The estimation gain for Ω_y is more substantial, especially for higher values of q . This is corroborated by our simulation outputs (Tables 1 and 2), where the joint estimates perform better for both sets of parameters, but the differences between RF errors obtained from joint and separate estimates tend to be lower for $\hat{\Omega}_y$ than $\hat{\beta}$.

5.2 Remaining Challenges

Our proposed framework for inference in complex multilayer networks still presents a couple of challenges, which can benefit from theoretical work pursued in the following directions. First, the use of tuning parameter selection criteria requires a rigorous analysis, not only to select parameter estimates that exhibit good finite sample performance, but also to ensure that the selected estimates when plugged into the hypothesis testing procedures result in tests and confidence intervals with accurate size, power or coverage guarantees. A number of existing methods have given consistency results of the (extended) BIC based tuning parameter selection procedures in high-dimensional regression (??) or graphical model (??) setups. Using theoretical tools provided therein and generalizing or adapting their conditions to multilayer settings is a possible avenue that can be explored.

To maintain focus on the current problem, we defer this to future work. Secondly, extending JMMLE to include overlapping within- or between-layer groups is of interest to tackle practical situations like multiple pathways sharing a number of common genes. The current algorithm involves fitting multiple group lasso models (using the R package `grpreg`), which can be replaced by alternative methods that can handle overlapping groups.

5.3 Extensions

There are two immediate extensions of our hypothesis testing framework.

(I) In recent work, ? proposed a framework to test for structural similarities and differences across multiple *single layer* GGMs. For K GGMs with precision matrices $\Omega^k = (\omega_{ii'}^k)_{i,i' \in \mathcal{I}_p}$, a 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 is developed, 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 (2.12)) can be used instead, to perform testing on the between-layer edges. Theoretical properties of this procedure can be derived using results in ?, 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 $K > 2$, detection of the following sets of inter-layer edges can be scientifically significant:

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

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 involved approach.

The strength of our proposed debiased estimator (3.1) is that only generic estimates of relevant model parameters that satisfy general rate conditions are necessary to obtain a valid asymptotic distribution. This translates to a high degree of flexibility in choosing the method of estimation. Our formulation based on sparsity assumptions (Section 2.2) is a specific way (motivated by applications in Omics data integration) 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 components in such multi-layer setups, low-rank, group-sparse or sparse methods (or a combination thereof) can be plugged into our alternating algorithm. Results analogous 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.

Finally, extending our framework to non-Gaussian data and graph Laplacian structures is of interest. As seen for the $K = 1$ case in ?, their alternating block algorithm continues to give comparable results under shrunken or truncated empirical distributions of Gaussian errors. Similar results may be possible in the general case, and improvements can come from modifying different parts of the estimation algorithm. For example, the estimation of the precision matrices based on restricted support sets using log-likelihoods in (2.6) can be replaced by methods like nonparanormal estimation (?) or regularized score matching (?). For graph Laplacian structures, generalization of recent work on multilayer models (??) in the lines of the JMMLE framework may be explored.

Appendix

A Proofs of main results

Proof of Theorem 2.3. The theorem is a generalization of Theorem 1 in ?. The proof follows directly from the proof of that theorem, substituting $(\widehat{B}^{(k)}, \widehat{\Theta}_\epsilon^{(k)}), (B^*, \Theta_\epsilon^*)$ and $(B^\infty, \Theta_\epsilon^\infty)$ therein with $(\widehat{\mathcal{B}}^{(t)}, \widehat{\Theta}_y^{(t)}), (\mathcal{B}_0, \Theta_{y0})$ and $(\mathcal{B}^\infty, \Theta_y^\infty)$, respectively, and their corresponding variations as required. \square

We use the following condition extensively while deriving the results that follow.

Condition 3 (Restricted eigenvalues). A symmetric matrix $\mathbf{M} \in \mathbb{M}(b, b)$ is said to satisfy the restricted eigenvalue or RE condition with parameters $\psi, \phi > 0$, denoted as curvature and tolerance, respectively, if

$$\boldsymbol{\theta}^T \mathbf{M} \boldsymbol{\theta} \geq \psi \|\boldsymbol{\theta}\|^2 - \phi \|\boldsymbol{\theta}\|_1^2$$

for all $\boldsymbol{\theta} \in \mathbb{R}^b$. In short, this is denoted by $\mathbf{M} \sim RE(\psi, \phi)$.

Starting from ?, different versions of the RE conditions have been proposed and used in high-dimensional analysis (????) to ensure that a covariance matrix satisfies a somewhat relaxed positive-definiteness condition.

Proof of Theorem 2.3. The proof strategy is as follows. We first show that given fixed $(\mathcal{X}, \mathcal{E})$, and some conditions on $\widetilde{\mathbf{E}}^k := \mathbf{Y}^k - \mathbf{X}^k \widetilde{\mathbf{B}}^k, k \in \mathcal{I}_K$, the bounds in Theorem 2.3 hold. We then show that for random $(\mathcal{X}, \mathcal{E})$, those conditions hold with probability approaching 1.

Assume fixed \mathcal{X}, \mathcal{E} and deterministic $\widetilde{\mathcal{B}} = \{\widetilde{\mathbf{B}}^k\}$, and the following conditions.

(A1) For $k \in \mathcal{I}_K$,

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

with C_β is non-negative and depends on \mathcal{B}_0 only.

(A2) For all $j \in \mathcal{I}_q$,

$$\frac{1}{n} \left\| (\widetilde{\mathbf{E}}_{-j}^k)^T \widetilde{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_\infty \leq \mathbb{Q} \left(C_\beta, \Sigma_{x0}^k, \Sigma_{y0}^k \right) \sqrt{\frac{\log(pq)}{n}},$$

where $\mathbb{Q}(C_\beta, \Sigma_{x0}^k, \Sigma_{y0}^k)$ is non-negative and depends on $\mathcal{B}_0, \Sigma_{x0}^k$ and Σ_{y0}^k only.

(A3) Denote $\widetilde{\mathbf{S}}^k = (\widetilde{\mathbf{E}}^k)^T \widetilde{\mathbf{E}}^k / n$. Then $\widehat{\mathbf{S}}^k \sim RE(\psi^k, \phi^k)$ with $Kq\phi \leq \psi/2$ where $\psi = \min_k \psi^k, \phi = \max_k \phi^k$.

Then the following hold

(I) Given the choice of tuning parameter

$$\gamma_n \geq 4\sqrt{|g_{\max}|} \mathbb{Q}_0 \sqrt{\frac{\log(pq)}{n}}; \quad \mathbb{Q}_0 := \max_{k \in \mathcal{I}_K} \mathbb{Q} \left(C_\beta, \Sigma_{x0}^k, \Sigma_{y0}^k \right)$$

$$\|\widehat{\Theta}_j - \Theta_{0,j}\|_F \leq 12\sqrt{s_j}\gamma_n/\psi, \quad (\text{A.1})$$

$$\sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\widehat{\theta}_{jj'}^{[g]} - \theta_{0,jj'}^{[g]}\| \leq 48s_j\gamma_n/\psi. \quad (\text{A.2})$$

$$|\text{supp}(\widehat{\Theta}_j)| \leq 128s_j/\psi \quad (\text{A.3})$$

(II) For the choice of tuning parameter $\gamma_n = 4\sqrt{|g_{\max}|}\mathbb{Q}_0\sqrt{\log(pq)/n}$,

$$\frac{1}{K} \sum_{k=1}^K \|\widehat{\Omega}_y^k - \Omega_{y0}^k\|_F \leq O\left(\mathbb{Q}_0 \sqrt{\frac{|g_{\max}|S}{K}} \sqrt{\frac{\log(pq)}{n}}\right) \quad (\text{A.4})$$

Condition (A1) holds by assumption. When \mathcal{X} and \mathcal{E} are random, the following proposition ensures that (A2) and (A3) hold with probabilities approaching to 1.

Consider deterministic $\widetilde{\mathcal{B}}$ satisfying assumption (A1), and conditions (E1), (E2) from the main paper. Then for sample size $n \gtrsim \log(pq)$ and $k \in \mathcal{I}_K$,

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

$$\psi^k = \frac{\Lambda_{\min}(\Sigma_{x0}^k)}{2}; \quad \phi^k = \frac{\psi^k \log p}{n} + 2C_\beta c_2 [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} \frac{\log(pq)}{n}$$

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 $j \in \mathcal{I}_q$

$$\left\| \frac{1}{n} (\widetilde{\mathbf{E}}_{-j}^k)^T \widetilde{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_\infty \leq \mathbb{Q}\left(C_\beta, \Sigma_{x0}^k, \Sigma_{y0}^k\right) \sqrt{\frac{\log(pq)}{n}}$$

with probability $\geq 1 - 1/p^{\tau_1 - 2} - 12c_1 \exp[-(c_2^2 - 1) \log(pq)] - 6c_4 \exp[-(c_5^2 - 1) \log(pq)]$, $c_4 > 0$, $c_5 > 1$, $\tau_1 > 2$, where

$$\begin{aligned} \mathbb{Q}\left(C_\beta, \Sigma_{x0}^k, \Sigma_{y0}^k\right) &= \left[2C_\beta^2 V_x^k + 4C_\beta c_2 [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2}\right] \sqrt{\frac{\log(pq)}{n}} + \\ &\quad c_5 \left[\Lambda_{\max}(\Sigma_{y0}^k) \sigma_{y0,j,-j}^k\right]^{1/2} \sqrt{\frac{\log q}{\log(pq)}} \end{aligned}$$

with $\sigma_{y0,j,-j}^k = \text{Var}(E_j - \mathbb{E}_{-j} \theta_{0,j})$, and

$$V_x^k = \sqrt{\frac{\log 4 + \tau_1 \log p}{c_x^k n}}; \quad c_x^k = \left[128(1 + 4\Lambda_{\max}(\Sigma_{x0}^k))^2 \max_i (\sigma_{x0,ii}^k)^2\right]^{-1}$$

We prove the main theorem by putting together Lemma A and Lemma A, and simplifying the constants $c_1 := 12c_1$, $c_4 := 6c_4$. \square

Proof of Theorem 2.3. The strategy is the same as in Theorem 2.3. We first establish the theorem statements hold for fixed \mathcal{X}, \mathcal{E} in the presence of certain regularity conditions, and then show that those conditions are satisfied with probability approaching 1 when \mathcal{X} and \mathcal{E} are random.

Assume fixed $(\mathcal{X}, \mathcal{E})$, and deterministic $\tilde{\Theta} = \{\tilde{\Theta}_j\}$, so that

(B1) For $j \in \mathcal{I}_q$,

$$\|\tilde{\Theta}_j - \Theta_{0,j}\|_F \leq C_\Theta \sqrt{\frac{\log q}{n}},$$

for some C_Θ dependent on Θ_0 only.

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

$$\left\| \tilde{\gamma}^k - \tilde{\Gamma}^k \beta_0 \right\|_\infty \leq \mathbb{R}(C_\Theta, \Sigma_{x0}^k, \Sigma_{y0}^k) \sqrt{\frac{\log(pq)}{n}}.$$

where $\mathbb{R}(C_\Theta, \Sigma_{x0}^k, \Sigma_{y0}^k)$ depends on Θ_0, Σ_{x0}^k and Σ_{y0}^k only, and $\{\tilde{\mathbf{T}}^k\}$ are defined using $\tilde{\Theta}$ similar to (2.14).

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

Then, given the choice of the tuning parameter

$$\lambda_n \geq 4\sqrt{|h_{\max}|} \mathbb{R}_0 \sqrt{\frac{\log(pq)}{n}}; \quad \mathbb{R}_0 := \max_{k \in \mathcal{I}_K} \mathbb{R}(C_\Theta, \Sigma_{x0}^k, \Sigma_{y0}^k)$$

the following holds

$$\|\hat{\beta} - \beta_0\|_1 \leq 48\sqrt{|h_{\max}|} B \lambda_n / \psi^* \quad (\text{A.5})$$

$$\|\hat{\beta} - \beta_0\| \leq 12\sqrt{B} \lambda_n / \psi^* \quad (\text{A.6})$$

$$\sum_{h \in \mathcal{H}} \|\beta^{[h]} - \beta_0^{[h]}\| \leq 48B \lambda_n / \psi^* \quad (\text{A.7})$$

$$(\hat{\beta} - \beta_0)^T \hat{\Gamma} (\hat{\beta} - \beta_0) \leq 72B \lambda_n^2 / \psi^* \quad (\text{A.8})$$

Condition (B1) holds by assumption. Next, we verify that conditions (B2) and (B3) hold with high probability given fixed $\tilde{\Theta}$.

Consider deterministic $\tilde{\Theta}$ satisfying assumption (B1). Also assume conditions (E3), (E4) from the main body of the paper. Then, for sample size $n \gtrsim \log(pq)$,

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

$$\psi_* = \min_k \psi^k \left(\min_i \psi_k^j - d_k C_\Theta \sqrt{\frac{\log(pq)}{n}} \right), \quad \phi_* = \max_k \phi^k \left(\min_i \phi_k^j + d_k C_\Theta \sqrt{\frac{\log(pq)}{n}} \right)$$

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

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

$$\mathbb{R}(C_\Theta, \Sigma_{x0}^k, \Sigma_{y0}^k) = c_2 \left\{ d_k C_\Theta \sqrt{\frac{\log(pq)}{n}} [\Lambda_{\max}(\Sigma_{x0}^k) \Lambda_{\max}(\Sigma_{y0}^k)]^{1/2} + \left[\frac{\Lambda_{\max}(\Sigma_{x0}^k)}{\Lambda_{\min}(\Sigma_{y0}^k)} \right]^{1/2} \right\}$$

The theorem follows straightforwardly by putting together the results from Lemmas A and A, and simplifying the constants $c_1 := 12c_1, c_4 := 6c_4$. \square

Proof of Theorem 2.3. The first part is immediate from the proof of part I of Theorem 4 in ?. By choice of λ_n , we now have

$$\|\widehat{\mathbf{B}}^{k(0)} - \mathbf{B}_0^k\|_1 = O\left(\sqrt{\frac{\log(pq)}{n}}\right),$$

so we can apply Theorem 2.3 to prove the bounds on $\{\widehat{\Theta}_j^{(0)}\}$. \square

Proof of Theorem 3.1. Define the following:

$$\widehat{\mathbf{D}}_i = \text{vec}(\widehat{\mathbf{b}}_i^1, \dots, \widehat{\mathbf{b}}_i^K); \quad \mathbf{R}_i^k = \mathbf{X}_i^k - \mathbf{X}_{-i}^k \widehat{\boldsymbol{\zeta}}_i^k; k \in \mathcal{I}_K$$

Then, from (3.1) we have

$$\mathbf{M}_i(\widehat{\mathbf{C}}_i - \widehat{\mathbf{D}}_i)^T = \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \widehat{\mathbf{E}}^1 \\ \vdots \\ \frac{1}{\widehat{s}_i^K} (\mathbf{R}_i^K)^T \widehat{\mathbf{E}}^K \end{bmatrix} \quad (\text{A.9})$$

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

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

Putting them back in (A.9) and using $t_i^k = (\mathbf{R}_i^k)^T \mathbf{X}_i^k / n$, we get

$$\begin{aligned} \mathbf{M}_i(\widehat{\mathbf{C}}_i - \widehat{\mathbf{D}}_i)^T &= \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{E}^1 \\ \vdots \\ \frac{1}{\widehat{s}_i^K} (\mathbf{R}_i^K)^T \mathbf{E}^K \end{bmatrix} + \mathbf{M}_i(\mathbf{D}_i - \widehat{\mathbf{D}}_i)^T \\ &\quad + \frac{1}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{X}_{-i}^1 (\mathbf{B}_{0,-i}^1 - \widehat{\mathbf{B}}_{-i}^1) \\ \vdots \\ \frac{1}{\widehat{s}_i^K} (\mathbf{R}_i^K)^T \mathbf{X}_{-i}^K (\mathbf{B}_{0,-i}^K - \widehat{\mathbf{B}}_{-i}^K) \end{bmatrix} \\ \Rightarrow \widehat{\Omega}_y^{1/2} \mathbf{M}_i(\widehat{\mathbf{C}}_i - \mathbf{D}_i)^T &= \frac{\widehat{\Omega}_y^{1/2}}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{E}^1 \\ \vdots \\ \frac{1}{\widehat{s}_i^K} (\mathbf{R}_i^K)^T \mathbf{E}^K \end{bmatrix} + \frac{\widehat{\Omega}_y^{1/2}}{\sqrt{n}} \begin{bmatrix} \frac{1}{\widehat{s}_i^1} (\mathbf{R}_i^1)^T \mathbf{X}_{-i}^1 (\mathbf{B}_{0,-i}^1 - \widehat{\mathbf{B}}_{-i}^1) \\ \vdots \\ \frac{1}{\widehat{s}_i^K} (\mathbf{R}_i^K)^T \mathbf{X}_{-i}^K (\mathbf{B}_{0,-i}^K - \widehat{\mathbf{B}}_{-i}^K) \end{bmatrix} \quad (\text{A.10}) \end{aligned}$$

At this point, we drop k and 0 in the subscripts since there is no ambiguity, and establish the following:

Given conditions (T1) and (T2), the following holds for sample size n such that $n \gtrsim \log(pq)$:

$$\frac{1}{\sqrt{n\hat{s}_i}} \hat{\Omega}_y^{1/2} \mathbf{E}^T \mathbf{R}_i \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{1n};$$

$$\|\mathbf{S}_{1n}\|_\infty \leq \frac{D_\Omega^{1/2}(2 + D_\zeta)c_2[\Lambda_{\max}(\Sigma_x)\Lambda_{\max}(\Sigma_e)]^{1/2}\sqrt{\log(pq)}}{\sqrt{\sigma_{x,i,-i}} - n^{-1/4} - D_\zeta\sqrt{V_x}} = O\left(\frac{\log(pq)}{\sqrt{n}}\right) \quad (\text{A.11})$$

with probability $\geq 1 - 6c_1 \exp[-(c_2^2 - 1)\log(pq)] - 1/p^{\tau_1-2} - \kappa_i/\sqrt{n}$, where $\kappa_i := \text{Var}[(X_i - \mathbb{X}_{-i}\boldsymbol{\zeta}_{0,-i})^2]$.

Additionally, given condition (T3) we have

$$\left\| \frac{1}{\sqrt{n\hat{s}_i}} \mathbf{R}_i^T \mathbf{X}_{-i} (\mathbf{B}_{-i} - \hat{\mathbf{B}}_{-i}) \hat{\Omega}_y^{1/2} \right\|_\infty$$

$$\leq \frac{D_\beta(\Lambda_{\min}(\Sigma_y)^{1/2} + D_\Omega^{1/2})}{\sigma_{x,i,-i} - n^{-1/2} - D_\zeta\sqrt{V_x}} \left[c_7 \sqrt{(\sqrt{\sigma_{x,i,-i}}\Lambda_{\max}(\Sigma_{x,-i}))\log p} + \sqrt{n}D_\zeta V_x \right] = O\left(\frac{\log(pq)}{\sqrt{n}}\right) \quad (\text{A.12})$$

holds with probability $\geq 1 - 6c_6 \exp[-(c_7^2 - 1)\log(pq)] - 1/p^{\tau_1-2} - \kappa_i/\sqrt{n}$ for some $c_6 > 0, c_7 > 1$.

Given Lemma A, the first and second summands on the right hand side of (A.10) are bounded above by applying each of (A.11) and (A.12) K times. This completes the proof. \square

Proof of Theorem 3.2. From (A.10) and Lemma A we have that

$$(\hat{\Omega}_y^k)^{1/2} m_i^k (\hat{\mathbf{c}}_i^k - \mathbf{b}_{0i}^k) \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{2n}^k, \quad (\text{A.13})$$

where $\|\mathbf{S}_{2n}^k\|_\infty = o_P(1)$. We next obtain the following lemma: Drop k in superscripts and 0 in subscripts. Given condition (T1), the following holds with probability $\geq 1 - 6c_6 \exp[-(c_7^2 - 1)\log(p-1)] - 1/p^{\tau_2-2} - \kappa_i/\sqrt{n}, \tau_2 > 2$:

$$\left| \frac{m_i}{\sqrt{n}} - \sqrt{\sigma_{x,i,-i}} \right| \leq \delta_i := \sqrt{\frac{\log 4 + \tau_2}{c_i n}} + \frac{D_\zeta + 1}{\sqrt{\sigma_{x,i,-i}} - n^{-1/2} - D_\zeta\sqrt{V_x}} \times$$

$$\left[c_7 [(\sigma_{x,i,-i}\Lambda_{\max}(\Sigma_{x,-i}))^{1/2} \sqrt{\frac{\log p}{n}} + D_\zeta V_x \right], \quad (\text{A.14})$$

where $c_i = [128(1 + 4\sigma_{x,i,-i})^2(\sigma_{x,i,-i})^2]^{-1}$, and the sample size satisfies $n \gtrsim \log p$.

We also provide the following general result:

Consider two positive definite matrices $\mathbf{A}, \mathbf{A}_1 \in \mathbb{M}(a, a)$. Then, for $\delta > 0$, we have

$$\|\mathbf{A} - \mathbf{A}_1\|_\infty \leq \delta \Rightarrow \|\mathbf{A}^{1/2} - \mathbf{A}_1^{1/2}\|_\infty \leq \delta^{1/2}.$$

Applying Lemma A it follows immediately from assumption (T2) that

$$\left\| \widehat{\Omega}_y^{1/2} - \Omega_y^{1/2} \right\|_{\infty} \leq D_{\Omega}^{1/2} \quad (\text{A.15})$$

Using Lemma A in conjunction with (A.15) we now have

$$\begin{aligned} \sqrt{n}(\Omega_{y0}^k)^{1/2} \sqrt{\sigma_{x0,i,-i}^k} (\widehat{\mathbf{c}}_i^k - \mathbf{b}_{0i}^k) &\sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{3n}^k \\ \Rightarrow \sqrt{n}\Sigma_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2 - \boldsymbol{\delta}) &\sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{3n}, \end{aligned} \quad (\text{A.16})$$

where $\Sigma_i := \Sigma_{y0}^1/\sigma_{x0,i,-i}^1 + \Sigma_{y0}^2/\sigma_{x0,i,-i}^2$ and $\mathbf{S}_{3n} = \mathbf{S}_{3n}^1 - \mathbf{S}_{3n}^2$, $\|\mathbf{S}_{3n}^k\|_{\infty} = o_P(1)$. We now break down the left hand side above as

$$\begin{aligned} \sqrt{n}\Sigma_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2 - \boldsymbol{\delta}) &= \sqrt{n}\Sigma_i^{-1/2}\widehat{\Sigma}_i^{1/2}\widehat{\Sigma}_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2) - \sqrt{n}\Sigma_i^{-1/2}\boldsymbol{\delta} \\ &= (\Sigma_i^{-1/2}\widehat{\Sigma}_i^{1/2} - \mathbf{I}) \cdot \sqrt{n}\widehat{\Sigma}_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2) + \\ &\quad \sqrt{n}\widehat{\Sigma}_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2) - \sqrt{n}\Sigma_i^{-1/2}\boldsymbol{\delta}, \end{aligned} \quad (\text{A.17})$$

with

$$\widehat{\Sigma}_i := \frac{n\widehat{\Sigma}_y^1}{(m_i^1)^2} + \frac{n\widehat{\Sigma}_y^2}{(m_i^2)^2}.$$

Next, we obtain the following lemma: Given conditions (T1) and (T2), for the pooled covariance matrix estimate $\widehat{\Sigma}_i$, we have

$$\left\| \widehat{\Sigma}_i - \Sigma_i \right\|_{\infty} = o(1),$$

for sample size $n \gtrsim \log p$. Lemma A now implies that $\|\widehat{\Sigma}_i^{1/2} - \Sigma_i^{1/2}\|_{\infty} = o(1)$. Putting this in the first summand of (A.17), then using (A.16) we get

$$\sqrt{n}\widehat{\Sigma}_i^{-1/2}(\widehat{\mathbf{c}}_i^1 - \widehat{\mathbf{c}}_i^2) - \sqrt{n}\Sigma_i^{-1/2}\boldsymbol{\delta} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I}) + \mathbf{S}_{4n},$$

with $\|\mathbf{S}_{4n}\|_{\infty} = o_P(1)$. The power of the global test follows as a consequence. Finally, the lower bound on the order of $\|\boldsymbol{\delta}\|$ holds because $n\boldsymbol{\delta}^T \Sigma_i^{-1} \boldsymbol{\delta} \geq n\|\boldsymbol{\delta}\|^2 \Lambda_{\min}(\Sigma_i^{-1})$, and

$$\Lambda_{\min}(\Sigma_i^{-1}) = \frac{\Lambda_{\max}(\Sigma_{y0}^1)}{\sigma_{x,i,-i}^1} + \frac{\Lambda_{\max}(\Sigma_{y0}^2)}{\sigma_{x,i,-i}^2}.$$

□

Proof of Theorem 3.3. The proof follows the general structure of Theorem 4.1 in ?, with two modifications. Firstly, we replace the bound in equation (12) of ? by a new deviation bound

$$P\left(\left|d_{ij} - \frac{\mu_j}{\sigma_j}\right| \geq t\right) = (1 - \Phi(t))(1 + o(1))$$

for any t , since $(d_{ij} - \mu_j)/\sigma_j \sim N(0, 1) + o_P(1)$ from Corollary 3.2. We replace $G_{\kappa}(t)$ in all following calculations in ? with $1 - \Phi(t)$. Secondly, we need to ensure that given

both Σ_{y0}^1 and Σ_{y0}^2 satisfy the condition (D1) or (D1*), the pooled covariance matrix $\Sigma_{y0}^1/\sigma_{x0,i,-i}^1 + \Sigma_{y0}^2/\sigma_{x0,i,-i}^2$ also does so.

For this, denote $c_k = \sigma_{x0,i,-i}^k$, $k = 1, 2$. Notice that for any $C_1, C_2 > 0$,

$$\begin{aligned} r_{jj'}^k \geq C_k &\Rightarrow \sigma_{y0,jj'}^k \geq (\sigma_{y0,jj}^k \sigma_{y0,j'j'}^k)^{1/2} C_k \\ &\Rightarrow \frac{\sigma_{y0,jj'}^1}{c_1} + \frac{\sigma_{y0,jj'}^2}{c_2} \geq \frac{(\sigma_{y0,jj}^1 \sigma_{y0,j'j'}^1)^{1/2} C_1}{c_1} + \frac{(\sigma_{y0,jj}^2 \sigma_{y0,j'j'}^2)^{1/2} C_2}{c_2} \\ &\Rightarrow \frac{\sigma_{y0,jj'}^1/c_1 + \sigma_{y0,jj'}^2/c_2}{(\sigma_{y0,jj}^1 \sigma_{y0,j'j'}^1)^{1/2}/c_1 + (\sigma_{y0,jj}^2 \sigma_{y0,j'j'}^2)^{1/2}/c_2} \geq \min\{C_1, C_2\}. \end{aligned}$$

It now follows that (D1) or (D1*) holds for the pooled covariance matrices. \square

B Proofs of auxiliary results

Proof of Lemma A. The proof has the same structure as the proof of Theorem 1 in ?, where consistency of the (single layer) JSEM estimates are established. Part (I) is analogous to part A.1 therein, but the proof strategy is completely different, which we provide in detail next. Our part (II) follows along similar lines as parts A.2 and A.3, incorporating the updated quantities from the first part (A.1). For this part of the proof, we provide an outline and leave details to the reader.

Proof of part (I). In its reparametrized version, (2.12) 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^{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

$$\begin{aligned} \sum_{k=1}^K \left| (\mathbf{d}_j^k)^T \widehat{\mathbf{S}}^k \mathbf{T}_{0,j}^k \right| &= \sum_{k=1}^K \left| \mathbf{d}_j^k \cdot \frac{1}{n} (\widehat{\mathbf{E}}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right| \\ &\leq \sum_{k=1}^K \left\| \frac{1}{n} (\widehat{\mathbf{E}}_{-j}^k)^T \widehat{\mathbf{E}}^k \mathbf{T}_{0,j}^k \right\|_{\infty} \|\mathbf{d}_{-j}^k\|_1 \\ &\leq \left[\sum_{j \neq j', g \in \mathcal{G}_y^{jj'}} \|\mathbf{d}_{jj'}^{[g]}\| \right] \mathbb{Q}_0 \sqrt{|g_{\max}|} \sqrt{\frac{\log(pq)}{n}} \end{aligned}$$

by assumption (A2). 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 γ_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 ≥ 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 $\mathbf{D}_j = (\mathbf{d}_j^1, \dots, \mathbf{d}_j^K)$. 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 (A3).

Combining the above with (B.3), we finally have

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

Since

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

The bounds in (A.1) and (A.2) are obtained by replacing the corresponding elements in (B.4).

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

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

using the KKT condition for (2.12) and assumption (A2). The choice of γ_n now ensures that the right hand side is $\geq 3|g|\gamma_n/4$. Hence,

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

using (B.3) and (B.4).

Proof of part (II). We denote the selected edge set for the k^{th} Y-network by \hat{E}^k . Denote its population version by E_0^k . Further, let

$$\tilde{\Omega}_y^k = \text{diag}(\Omega_{y0}^k) + \Omega_{y, E_0^k \cap \hat{E}^k}^k$$

Based on similar derivations as in the proof of Corollary A.1 in ?, the following two upper bounds can be established:

$$|\hat{E}^k| \leq \frac{128S}{\psi} \quad (\text{B.5})$$

$$\frac{1}{K} \sum_{k=1}^K \|\tilde{\Omega}_y^k - \Omega_{y0}^k\|_F \leq \frac{12c_y \sqrt{S} \gamma_n}{\sqrt{K} \psi} \quad (\text{B.6})$$

following which, taking $\gamma_n = 4\sqrt{|g_{\max}|}\mathbb{Q}_0\sqrt{\log(pq)/n}$,

$$\Lambda_{\min}(\tilde{\Omega}_y^k) \geq d_y - \frac{48c_y\mathbb{Q}_0\sqrt{|g_{\max}|}S}{\psi}\sqrt{\frac{\log(pq)}{n}} \geq (1-t_1)d_y > 0 \quad (\text{B.7})$$

$$\Lambda_{\max}(\tilde{\Omega}_y^k) \leq c_y + \frac{48c_y\mathbb{Q}_0\sqrt{|g_{\max}|}S}{\psi}\sqrt{\frac{\log(pq)}{n}} \leq c_y + t_1d_y < \infty \quad (\text{B.8})$$

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

$$n \geq |g_{\max}|S \left[\frac{48c_y\mathbb{Q}_0}{\psi t_1 d_y} \right]^2 \log(pq).$$

Following the same steps as part A.3 in the proof of Theorem 4.1 in ?, it can be proven using (B.5)–(B.8) that

$$\sum_{k=1}^K \left\| \hat{\Omega}_y^k - \tilde{\Omega}_y^k \right\|_F^2 \leq O \left(\mathbb{Q}_0^2 |g_{\max}| S \frac{\log(pq)}{n} \right)$$

The proof is now complete by combining this with (B.6) and then applying the Cauchy-Schwarz inequality and the triangle inequality. \square

Proof of Lemma A. We drop the subscript 0 for true values and the superscript k since there is no scope of ambiguity. For part 1, we start with an auxiliary lemma: For a sub-Gaussian design matrix $\mathbf{X} \in \mathbb{M}(n, p)$ with columns having mean $\mathbf{0}_p$ and covariance matrix Σ_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 $c_3 > 0$. Denote $\hat{\mathbf{E}} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}$. For $\mathbf{v} \in \mathbb{R}^q$, we have

$$\begin{aligned} \mathbf{v}^T \hat{\mathbf{S}} \mathbf{v} &= \frac{1}{n} \|\hat{\mathbf{E}} \mathbf{v}\|^2 \\ &= \frac{1}{n} \|(\mathbf{E} + \mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}})) \mathbf{v}\|^2 \\ &= \mathbf{v}^T \mathbf{S} \mathbf{v} + \frac{1}{n} \|\mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}}) \mathbf{v}\|^2 + 2\mathbf{v}^T (\mathbf{B}_0 - \hat{\mathbf{B}})^T \left(\frac{(\mathbf{X})^T \mathbf{E}}{n} \right) \mathbf{v} \end{aligned} \quad (\text{B.9})$$

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

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

by assumption (A1). Now, we use another lemma: For zero-mean independent sub-gaussian matrices $\mathbf{X} \in \mathbb{M}(n, p)$, $\mathbf{E} \in \mathbb{M}(n, q)$ with parameters (Σ_x, σ_x^2) and (Σ_e, σ_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 (B.9) and get

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

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:

$$\begin{aligned} \left\| \frac{1}{n} \hat{\mathbf{E}}_{-j}^T \hat{\mathbf{E}} \mathbf{T}_{0,j} \right\|_\infty &= \left\| \frac{1}{n} [\mathbf{E}_{-j} + \mathbf{X}(\mathbf{B}_{0,j} - \hat{\mathbf{B}}_j)]^T [\mathbf{E} + \mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}})] \mathbf{T}_{0,j} \right\|_\infty \\ &\leq \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{E} \mathbf{T}_{0,j} \right\|_\infty + \left\| \frac{1}{n} \mathbf{E}_{-j}^T \mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}}) \mathbf{T}_{0,j} \right\|_\infty \\ &\quad + \left\| \frac{1}{n} (\mathbf{B}_{0,j} - \hat{\mathbf{B}}_j)^T \mathbf{X}^T \mathbf{X}(\mathbf{B}_0 - \hat{\mathbf{B}}) \mathbf{T}_{0,j} \right\|_\infty + \left\| \frac{1}{n} (\mathbf{B}_{0,j} - \hat{\mathbf{B}}_j)^T \mathbf{X}^T \mathbf{E} \mathbf{T}_{0,j} \right\|_\infty \\ &= \|\mathbf{W}_1\|_\infty + \|\mathbf{W}_2\|_\infty + \|\mathbf{W}_3\|_\infty + \|\mathbf{W}_4\|_\infty \end{aligned} \quad (\text{B.10})$$

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 B on \mathbf{E}_{-j} and $\mathbf{E}_j - \mathbf{E}_{-j} \boldsymbol{\theta}_{0,j}$ to obtain that for $n \gtrsim \log(q-1)$,

$$\|\mathbf{W}_1\|_\infty \leq c_5 [\Lambda_{\max}(\Sigma_{y,-j}) \sigma_{y,j,-j}]^{1/2} \sqrt{\frac{\log(q-1)}{n}} \quad (\text{B.11})$$

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

For \mathbf{W}_2 and \mathbf{W}_4 , identical bounds hold:

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

Since Ω_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 $n \gtrsim \log(pq)$,

$$\|\mathbf{W}_2\|_\infty + \|\mathbf{W}_4\|_\infty \leq 4C_\beta c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_y)]^{1/2} \frac{\log(pq)}{n} \quad (\text{B.12})$$

with probability $\geq 1 - 12c_1 \exp[-(c_2^2 - 1) \log(pq)]$ by applying Lemma B and assumption (A1).

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

$$n \geq 512(1 + 4\Lambda_{\max}(\Sigma_x^k))^4 \max_i (\sigma_{x,ii}^k)^4 \log(4p^{\tau_1}) \quad (\text{B.13})$$

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,

$$\|\mathbf{W}_4\|_\infty \leq \left\| \frac{\mathbf{X}^T \mathbf{X}}{n} \right\|_\infty \|\hat{\mathbf{B}} - \mathbf{B}_0\|_1^2 \|\mathbf{T}_{0,j}\|_1 \leq 2C_\beta^2 V_x \frac{\log(pq)}{n} \quad (\text{B.14})$$

We now bound the right hand side of (B.10) using (B.11), (B.12) and (B.14) to complete the proof, with the leading term of the sample size requirement being $n \gtrsim \log(pq)$. \square

Proof of Lemma A. The proof follows that of part (I) of Lemma A, with a different group norm structure. We only point out the differences.

Putting $\beta = \beta_0$ in (2.13) we get

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

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

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

Proceeding similarly as the proof of part (I) of Lemma A, with a different deviation bound and choice of λ_n , we get expressions equivalent to (B.3) and (B.4) respectively:

$$\mathbf{b}^T \hat{\Gamma} \mathbf{b} \leq \frac{3}{2} \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\| \quad (\text{B.15})$$

$$\frac{\psi^*}{3} \|\mathbf{b}\|^2 \leq \lambda_n \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\| \leq 4\lambda_n \sqrt{B} \|\mathbf{b}\| \quad (\text{B.16})$$

Furthermore, $\|\mathbf{b}\|_1 \leq \sqrt{|h_{\max}|} \sum_{h \in \mathcal{H}} \|\mathbf{b}^{[h]}\|$. The bounds in (A.5), (A.6), (A.7) and (A.8) now follow. \square

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

$$\hat{\mathbf{T}}_k^2 \otimes \hat{\Sigma}_x^k \sim RE(\psi_*^k, \phi_*^k) \quad (\text{B.17})$$

with high enough probability. because then we can take $\psi_* = \min_k \psi_*^k, \phi_* = \max_k \phi_*^k$. The proof of (B.17) follows similar lines of the proof of Proposition 1 in ?, only replacing $\Theta_\epsilon, \hat{\Theta}_\epsilon, \mathbf{X}$ therein with $(\mathbf{T}^k)^2, (\hat{\mathbf{T}}^k)^2, \mathbf{X}^k$, respectively. We omit the details.

Part 2 follows the proof of Proposition 2 in ?. \square

Proof of Lemma A. To show (A.11) we have

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

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

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

because Ω_x is diagonally dominant implies $\left\| \zeta_i \right\|_1 = \sum_{i' \neq i} |\omega_{x,ii'}| / \omega_{x,ii} \leq 1$, and using assumption (T1) and (A.15). Applying Lemma B, the following holds for $n \gtrsim \log(pq)$:

$$\frac{1}{\sqrt{n}} \left\| (\hat{\Omega}_y^{1/2} - \Omega_y^{1/2}) \mathbf{E}^T \mathbf{R}_i \right\|_\infty \leq \sqrt{D_\Omega} (2 + D_\zeta) c_2 [\Lambda_{\max}(\Sigma_x) \Lambda_{\max}(\Sigma_e)]^{1/2} \sqrt{\log(pq)} \quad (\text{B.18})$$

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

On the other hand,

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

which implies $s_i \leq \hat{s}_i + D_\zeta \sqrt{V_x}$. By applying Lemma 8 of ?,

$$\left\| \frac{1}{n} \mathbf{X}_{-i}^T \mathbf{X}_{-i} \right\|_\infty \leq \left\| \frac{1}{n} \mathbf{X}^T \mathbf{X} \right\|_\infty \leq V_x \quad (\text{B.19})$$

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

$$n \geq 512(1 + 4\Lambda_{\max}(\Sigma_x))^4 \max_i (\sigma_{x,ii})^4 \log(4p^{\tau_1}) \quad (\text{B.20})$$

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

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

Taking $\epsilon = n^{-1/4}$, we have $s_i \geq \sqrt{\sigma_{x,i,-i}} - n^{-1/4}$ with probability $\geq 1 - \kappa_i n^{-1/2}$. Then, for n satisfying (B.20) and $\sqrt{\sigma_{x,i,-i}} - n^{-1/4} > D_\zeta \sqrt{V_x}$, we get the bound with the above probability:

$$\frac{1}{\widehat{s}_i} \leq \frac{1}{\sqrt{\sigma_{x,i,-i}} - n^{-1/4} - D_\zeta \sqrt{V_x}} \quad (\text{B.21})$$

Combining (B.18) and (B.21) gives the upper bound for the right hand side of (A.11) with the requisite probability and sample size conditions.

To prove (A.12) we have

$$\frac{1}{n} \|\mathbf{R}_i^T \mathbf{X}_{-i}\|_\infty \leq \frac{1}{n} \|(\mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_{0,i})^T \mathbf{X}_{-i}\|_\infty + \frac{1}{n} \|\mathbf{X}_{-i}^T \mathbf{X}_{-i} (\widehat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i})\|_\infty \quad (\text{B.22})$$

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

$$\frac{1}{n} \|(\mathbf{X}_i - \mathbf{X}_{-i} \boldsymbol{\zeta}_i)^T \mathbf{X}_{-i}\|_\infty \leq c_7 [\sigma_{x,i,-i} \Lambda_{\max}(\Sigma_{x,-i})]^{1/2} \sqrt{\frac{\log(p-1)}{n}} \quad (\text{B.23})$$

with probability $\geq 1 - 6c_6 \exp[-(c_7^2 - 1) \log(p-1)]$ for some $c_6 > 0, c_7 > 1$. By (B.19), the second term on the right side of (B.22) is bounded above by $D_\zeta V_x$ with probability $\geq 1 - 1/p^{\tau_1-2}$ and n satisfying (B.20). The bound of (A.12) now follows by conditions (T2), (T3) and (B.21). Since $\sqrt{\sigma_{x,i,-i}} - n^{-1/4} > D_\zeta \sqrt{V_x}$ implies $\sqrt{\sigma_{x,i,-i}} > D_\zeta \sqrt{V_x}$, and $D_\zeta = O(\sqrt{\log p/n})$, the leading term of the overall sample size requirement is $n \gtrsim \log(pq)$. \square

Proof of Lemma A. We drop k in the superscripts. By definition,

$$\begin{aligned} \frac{m_i}{\sqrt{n}} &= \frac{1}{\widehat{s}_i} \frac{(\mathbf{X}_i - \mathbf{X}_{-i} \widehat{\boldsymbol{\zeta}}_i)^T \mathbf{X}_i}{n} \\ &= \frac{1}{\widehat{s}_i} \left[\frac{\|\mathbf{X}_i - \mathbf{X}_{-i} \widehat{\boldsymbol{\zeta}}_i\|^2}{n} + \frac{(\mathbf{X}_i - \mathbf{X}_{-i} \widehat{\boldsymbol{\zeta}}_i)^T \mathbf{X}_{-i} \widehat{\boldsymbol{\zeta}}_i}{n} \right] \\ &\leq \widehat{s}_i + \frac{1}{\widehat{s}_i} \cdot \frac{1}{n} \|\mathbf{R}_i^T \mathbf{X}_{-i}\|_\infty \left(\|\widehat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_{0,i}\|_1 + \|\boldsymbol{\zeta}_{0,i}\|_1 \right) \\ \Rightarrow \left| \frac{m_i}{\sqrt{n}} - \sqrt{\sigma_{x,i,-i}} \right| &\leq |\widehat{s}_i - \sqrt{\sigma_{x,i,-i}}| + \frac{1}{\widehat{s}_i} \cdot \frac{1}{n} \|\mathbf{R}_i^T \mathbf{X}_{-i}\|_\infty \left(\|\widehat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_i\|_1 + \|\boldsymbol{\zeta}_i\|_1 \right) \quad (\text{B.24}) \end{aligned}$$

By applying Lemma 8 in ?, we have a bound for the first summand on the right hand side:

$$|\widehat{s}_i - \sqrt{\sigma_{x,i,-i}}| \leq \sqrt{\frac{\log 4 + \tau_2}{c_i n}}; \quad c_i = [128(1 + 4\sigma_{x,i,-i})^2 \sigma_{x,i,-i}^2]^{-1},$$

with probability $1 - 1/p^{\tau_2-2}$ for some $\tau_2 > 2$, and $n \geq 512(1 + 4\sigma_{x,i,-i})^4 2\sigma_{x,i,-i}^4 \log(4)$. For the second summand in the right-hand side of (B.24), $1/\widehat{s}_i$ can be bounded using

(B.21), $(1/n)\|\mathbf{R}_i^T \mathbf{X}_{-i}\|_\infty$ can be bounded using derivations following (B.22). Finally, $\|\hat{\boldsymbol{\zeta}}_i - \boldsymbol{\zeta}_i\|_1 \leq D_\zeta$ from assumption (T1), and $\|\boldsymbol{\zeta}_i\|_1 \leq 1$ because Ω_x is diagonally dominant and $|\zeta_{ii'}| = |\omega_{x,ii'}|/\omega_{x,ii}$ for $i' \neq i$. The lemma now follows by putting everything back together in (B.24). \square

Proof of Lemma A. $\|\mathbf{A} - \mathbf{A}_1\|_\infty \leq \delta$ implies that $\mathbf{A}_1 + \delta \mathbf{J}_a \geq \mathbf{A}$ and $\mathbf{A} + \delta \mathbf{J}_a \geq \mathbf{A}_1$, where $\mathbf{J}_a \in \mathbb{M}(a, a)$ has all entries 1, and for positive definite matrices \mathbf{P}, \mathbf{Q} , $\mathbf{P} \geq \mathbf{Q}$ means $\mathbf{P} - \mathbf{Q}$ is positive definite. Now applying Theorem 1 part (a) in ? we have

$$(\mathbf{A} + \delta \mathbf{J}_a)^{1/2} \geq \mathbf{A}_1^{1/2}; \quad (\mathbf{A}_1 + \delta \mathbf{J}_a)^{1/2} \geq \mathbf{A}^{1/2}.$$

Using the same result, it is easy to prove that

$$\mathbf{A}^{1/2} + \sqrt{\delta} \mathbf{J}_a \geq (\mathbf{A} + \delta \mathbf{J}_a)^{1/2},$$

and the same for \mathbf{A}_1 . The lemma follows. \square

Proof of Lemma A. We drop k in the superscripts and 0 in subscripts. Note that it is enough to prove

$$\left\| \frac{n\hat{\Sigma}_y}{(m_i)^2} - \frac{\Sigma_y}{\sigma_{x,i,-i}} \right\|_\infty = o_P(1).$$

For this, consider the decomposition

$$\begin{aligned} \frac{n\hat{\Sigma}_y}{(m_i)^2} &= \frac{\hat{\Sigma}_y - \Sigma_y + \Sigma_y}{\sigma_{x,i,-i}} \cdot \frac{\sigma_{x,i,-i}}{(m_i)^2/n} \\ \Rightarrow \frac{n\hat{\Sigma}_y}{(m_i)^2} - \frac{\Sigma_y}{\sigma_{x,i,-i}} &= \frac{\hat{\Sigma}_y - \Sigma_y}{(m_i)^2/n} + \frac{\Sigma_y}{\sigma_{x,i,-i}} \left[1 - \frac{\sigma_{x,i,-i}}{(m_i)^2/n} \right] \\ &= \frac{n}{(m_i)^2} \left[\hat{\Sigma}_y - \Sigma_y + \frac{\Sigma_y}{\sigma_{x,i,-i}} \left(\frac{(m_i)^2}{n} - \sigma_{x,i,-i} \right) \right]. \end{aligned}$$

From Lemma A we now have

$$\frac{m_i}{\sqrt{n}} \geq \sqrt{\sigma_{x,i,-i}} - \delta_i \quad \Rightarrow \quad \frac{m_i^2}{n} \geq (\sqrt{\sigma_{x,i,-i}} - \delta_i)^2 \geq \sigma_{x,i,-i} - \delta_i^2,$$

so that

$$\left\| \frac{n\hat{\Sigma}_y}{(m_i)^2} - \frac{\Sigma_y}{\sigma_{x,i,-i}} \right\|_\infty \leq \frac{\|\hat{\Sigma}_y - \Sigma_y\|_\infty + \sigma_{x,i,-i}^{-1} \delta_i^2 \|\Sigma_y\|_\infty}{\sigma_{x,i,-i} - \delta_i^2}, \quad (\text{B.25})$$

with probability $\geq 1 - 6c_6 \exp[-(c_7^2 - 1) \log(p-1)] - 1/p^{\tau_2-2} - \kappa_i/\sqrt{n}$ and for sample size satisfying $n \gtrsim \log p$, $n \geq 512(1+4\sigma_{x,i,-i})^4 2(\sigma_{x,i,-i})^4 \log(4)$ and $\sqrt{\sigma_{x,i,-i}} > \max\{\delta_i, n^{-1/4} - D_\zeta \sqrt{V_x}\}$. For the ℓ_∞ norms on the right-hand side, we have

$$\|\Sigma_y\|_\infty = \|\Omega_y^{-1}\|_\infty \leq (\Delta_0(\Omega_y))^{-1} \quad (\text{B.26})$$

following ?. For a bound on $\|\hat{\Sigma}_y - \Sigma_y\|_\infty$, if condition (II) of Theorem 3.2 is satisfied then we have

$$\|\hat{\Sigma}_y - \Sigma_y\|_\infty \leq \tilde{D}_\Omega \quad (\text{B.27})$$

where $\tilde{D}_\Omega = O(D_\Omega)$ and $D_\Omega = O(\tilde{D}_\Omega)$?. If condition (I) is satisfied, denote $\epsilon = D_\Omega/\Delta_0(\Omega_y)$. Then

$$\begin{aligned} \|\hat{\Sigma}_y - \Sigma_y\|_\infty &= \|\hat{\Sigma}_y(\Omega_y - \hat{\Omega}_y)\Sigma_y\|_\infty \\ &\leq \|\hat{\Sigma}_y\|_\infty \|\Omega_y - \hat{\Omega}_y\|_\infty \|\Sigma_y\|_\infty \\ &\leq \|(\mathbf{I} + (\Omega_y - \hat{\Omega}_y)\Sigma_y)^{-1}\|_\infty \|\Sigma_y\|_\infty \epsilon \\ &\leq \frac{\epsilon}{\Delta_0(\Omega_y)} \left[1 + \sum_{t=1}^{\infty} (\|(\Omega_y - \hat{\Omega}_y)\Sigma_y\|_\infty)^t \right] \\ &\leq \frac{\epsilon}{(1 - \epsilon)\Delta_0(\Omega_y)} \\ &= \frac{D_\Omega}{(\Delta_0(\Omega_y) - D_\Omega)\Delta_0(\Omega_y)} \end{aligned} \quad (\text{B.28})$$

Combining (B.26) with (B.27) or (B.28) as required and putting them back in the right-hand side of (B.25), we get the needed. \square

Proof of Lemma B. This is the same as in Lemma 2 in Appendix B of ? and its proof can be found there. \square

Proof of Lemma B. This is a part of Lemma 3 of Appendix B in ?, and is proved therein. \square