

# High-dimensional functional graphical model structure learning via neighborhood selection approach

Boxin Zhao<sup>1</sup>, Percy S. Zhai<sup>1</sup>, Y. Samuel Wang<sup>2</sup> and Mladen Kolar<sup>3</sup>

<sup>1</sup>*Booth School of Business, University of Chicago*  
e-mail: [boxinz@uchicago.edu](mailto:boxinz@uchicago.edu); [percy.zhai@uchicago.edu](mailto:percy.zhai@uchicago.edu)

<sup>2</sup>*Department of Statistics and Data Science, Cornell University*  
e-mail: [ysw7@cornell.edu](mailto:ysw7@cornell.edu)

<sup>3</sup>*Department of Data Sciences and Operations, Marshall School of Business, University of Southern California*  
e-mail: [mkolar@usc.edu](mailto:mkolar@usc.edu)

**Abstract:** Undirected graphical models are widely used to model the conditional independence structure of vector-valued data. However, in many modern applications, for example those involving EEG and fMRI data, observations are more appropriately modeled as multivariate random functions rather than vectors. Functional graphical models have been proposed to model the conditional independence structure of such functional data. We propose a neighborhood selection approach to estimate the structure of Gaussian functional graphical models, where we first estimate the neighborhood of each node via a function-on-function regression and subsequently recover the entire graph structure by combining the estimated neighborhoods. Our approach only requires assumptions on the conditional distributions of random functions, and we estimate the conditional independence structure directly. We thus circumvent the need for a well-defined precision operator that may not exist when the functions are infinite dimensional. Additionally, the neighborhood selection approach is computationally efficient and can be easily parallelized. The statistical consistency of the proposed method in the high-dimensional setting is supported by both theory and experimental results. In addition, we study the effect of the choice of the function basis used for dimensionality reduction in an intermediate step. We give a heuristic criterion for choosing a function basis and motivate two practically useful choices, which we justify by both theory and experiments.

**MSC2020 subject classifications:** Primary 62H22, 62J07; secondary 62P10.

**Keywords and phrases:** Functional graphical model, neighborhood selection, fMRI data.

Received October 2022.

## Contents

1	Introduction . . . . .	1043
	1.1 Related work . . . . .	1045

---

arXiv: [2105.02487](https://arxiv.org/abs/2105.02487)

1.2	Notation . . . . .	1047
1.3	Outline of the paper . . . . .	1048
2	Methodology . . . . .	1049
2.1	Functional graphical model . . . . .	1049
2.2	Functional neighborhood selection . . . . .	1050
2.3	Vector-on-vector regression . . . . .	1052
2.4	Choice of basis . . . . .	1054
2.5	Selection of tuning parameters . . . . .	1057
3	Optimization algorithm . . . . .	1058
4	Theoretical properties . . . . .	1060
4.1	Prior fixed function basis . . . . .	1062
4.2	Data-dependent function basis . . . . .	1066
4.3	Theoretical guidance on the choice of function basis . . . . .	1068
4.3.1	Minimize $\omega(M)$ . . . . .	1069
4.3.2	Minimize an approximate upper bound . . . . .	1069
5	Simulations . . . . .	1070
5.1	Comparison with baseline methods . . . . .	1072
5.2	The effect of $\epsilon_n$ . . . . .	1074
5.3	Performance of cross-validation . . . . .	1077
6	Data analysis . . . . .	1077
7	Conclusion . . . . .	1083
A	Technical proofs . . . . .	1084
A.1	Proof of Theorem 2.1 . . . . .	1084
A.2	Derivation of (10) and (15) . . . . .	1086
A.3	Simplification of ADMM optimization problems . . . . .	1087
A.4	Derivation of (40) . . . . .	1088
A.5	Proposition A.1 and its proof . . . . .	1090
A.6	Proof of Theorem 4.1 . . . . .	1091
A.7	Proof of Theorem 4.3 . . . . .	1095
B	Useful lemmas . . . . .	1101
C	Wall-clock runtime comparison . . . . .	1118
D	Labels of ROIs in the AAL atlas . . . . .	1120
E	Table of notations . . . . .	1122
	Acknowledgments . . . . .	1124
	Funding . . . . .	1124
	Supplementary Material . . . . .	1124
	References . . . . .	1124

## 1. Introduction

Multivariate functional data are collected in applications such as neuroscience, medical science, traffic monitoring, and finance. Although each observation is typically only recorded at a discrete set of time points, the underlying process may be interpreted as a realization of a multivariate stochastic process in continuous time. Such interpretation can provide a unifying approach to the analysis

of classical functional data and longitudinal data, where functional data can be used to deal with sparsely observed measurements with noise [13].

Our work is centered around elucidating the conditional independence structure of multivariate random functions. Gaining a robust understanding of such a structure can yield extensive applications, including the interpretation of time course gene expression data in genomics [66], multivariate time series data in finance [61], and electroencephalography (EEG) data in neuroscience [52, 53]. This paper is motivated by the analysis of data gathered from fMRI scans conducted on 116 distinct brain regions, with a time-series signal recorded for each region [45]. The sample comprises two groups: one group of individuals diagnosed with Attention Deficit Hyperactivity Disorder (ADHD) and a control group. Our aim is to comprehend the functional connectivity patterns between these brain regions for both the ADHD and control groups. Such functional connectivity can be uncovered by determining the conditional independence structure across the 116 random functions.

Graphical models are widely used to represent the conditional independence structure of multivariate random variables [32]. Let  $G = \{V, E\}$  denote an undirected graph where  $V$  is the set of vertices and  $E \subset V^2$  is the set of edges. When the data consist of random vectors  $\mathbf{X} = (X_1, \dots, X_p)^\top$ , we say that  $\mathbf{X}$  satisfies the pairwise Markov property with respect to  $G$  if  $X_v \perp\!\!\!\perp X_w \mid \{X_u\}_{u \in V \setminus \{v, w\}}$  holds if and only if  $\{v, w\} \notin E$ . This notion has been extended to *functional graphical models*—where each node represents a random function rather than a random scalar—in order to characterize the conditional independence relationship of multivariate random functions.

We propose a procedure to estimate the functional graphical model when the random functions follow a *multivariate Gaussian processes* (MGP). This setting was considered in [52], who proposed the functional graphical lasso to estimate the structure of the graph. Their procedure first obtains a finite dimensional representation for the observed multivariate functions using functional principal component analysis (FPCA). Subsequently, a precision matrix is computed from the projection scores of the finite dimensional representation using a graphical lasso objective with a group penalty. The graph structure is finally obtained from the non-zero blocks of the estimated precision matrix. When the underlying random functions are infinite dimensional, the corresponding covariance operator is a compact operator, and its inverse, the precision operator, is ill-defined [22]. As a result, [52] ensure their estimand is well defined by requiring that the random functions lie in a finite dimensional space. However, that assumption is restrictive and excludes infinite dimensional functional data.

In contrast to the functional graphical lasso proposed by [52], we propose a neighborhood selection approach to estimate Gaussian functional graphical models. For vector-valued Gaussian graphical models, [42] proposed a neighborhood selection procedure that estimates the neighborhood—the set of adjacent nodes in a conditional independence graph—for each node separately by sparse regression. The entire graph structure is then estimated by combining estimates of node-specific neighborhoods. We extend their approach to the functional data setting. This allows us to avoid defining the precision operator, and, as a result,

our theory extends to truly infinite dimensional functional data.

We cast the neighborhood selection procedure as a function-on-function regression problem. Due to the infinite dimensional nature of the functional data, we first project all observed random functions onto a finite dimensional basis. Thus, we approximate the function-on-function regression with a vector-on-vector regression problem that is solved by minimizing a squared error loss with a group lasso penalty. We do not require a specific choice of function basis for our methodology and the corresponding theory, and we provide a theoretically guided intuition on the choice of function basis under different conditions. Specifically, when estimating the neighborhood of a target node, we project all functions onto a single subspace instead of projecting each function onto its own subspace. In Section 2.4 we provide intuition for why this may be preferable to projecting each function onto its own subspace. This is supported by both the theory in Section 4.3 and the simulations in Section 5.

In addition to the methodology, we also provide nontrivial theoretical contributions. Most importantly, by directly estimating the conditional independence structure without reference to a population “precision operator,” we do not require that the functional data are finite dimensional. However, in the infinite dimensional setting, there will be a residual term due to using a finite dimensional approximation, and deriving error bounds requires a careful analysis of this residual term. Finally, our theory is non-asymptotic in nature, and we derive finite-sample guarantees for graph recovery.

In summary, the neighborhood selection approach yields at least three benefits. First, it allows us to define functional graphical models directly from the conditional distribution and does not require the notion of a precision operator. As a result, we can estimate the graph structure even from infinite dimensional data, rather than restricting the data to finite dimensional functions. Second, by estimating the neighborhood of each node separately, we have increased flexibility in choosing the function basis used to represent the random functions, and tailoring the function basis for the specific task at hand results in empirically better estimation results. Finally, when estimating the neighborhood of a node, we only need to handle  $p$  individual  $M \times M$  matrices. These neighborhood estimation procedures can be performed in parallel, leading to a highly efficient estimation procedure. In comparison, fglasso [52] needs to estimate a  $pM \times pM$  matrix, which is computationally much more expensive. We demonstrate the practical value of our neighborhood selection method on the motivating ADHD fMRI data set, and also on another ASD fMRI data set.

### 1.1. Related work

Our paper contributes to the growing literature on modeling multivariate functional data. We study the estimation of the conditional independence structure from multivariate functional data in the setting of MGPs [52]. For each component of the multivariate functional data, [52] projected observed functions on the corresponding function basis estimated by FPCA. Subsequently, the structure of the graph is estimated from the projection scores using the multitask

extension of the graphical lasso [28, 29], which estimates the precision matrix with a block structure. However, the precision operator is ill-defined when the functional data is infinite dimensional, and their method is computationally expensive when the number of nodes is large.

In the same setting, [53] proposed a dynamic functional graphical model that allows the graph structure to change over time. [72] proposed a Bayesian approach to functional graphical models. [67] studied estimation of the graph structure under the assumption of partial separability. Roughly speaking, partial separability assumes that the time-varying covariance of the MGP can be decomposed node-wise into a time-varying component and a constant. However, this assumption can be restrictive and may not hold in many settings; [40] proposed a test to verify the partial separability assumption. [68] and [69] discussed direct estimation of the difference between two Gaussian functional graphical models without the need to estimate each individual structure. [60] studied a latent multi-modal functional graphical model. [59] extended the Gaussian functional graphical model to a copula version by allowing monotonic transformations of the FPCA scores. In addition, [37] and [35] discussed a nonparametric functional graphical model; however, the graph therein is defined based on the additive conditional independence (ACI) relationship [36, 33, 34], which is not equivalent to the conditional independence (CI) relationship and is thus not directly comparable to our paper.

Our paper is also related to the literature on *function-on-function regression* that studies regression problems in which both the response and predictor variables are functions. [38] and [50] studied function-on-function linear regression, where each predictor function is transformed by a corresponding integral operator defined by the bivariate coefficient function, and the addition of all transformed predictor functions is defined as the signal function. The response function is then assumed to be a simple addition of the signal, intercept, and noise functions. To estimate the coefficient functions, they used the FPCA basis of the signal function to expand the observed functions and transform the original function-on-function regression problem into a function-on-scalar regression problem with uncorrelated predictors. Ultimately, this regression is solved by a penalized least-squares method. In contrast, we focus on variable selection, rather than prediction, and develop non-asymptotic theory. [39] also converted the function-on-function regression problem to function-on-scalar regression by projecting predictor functions with basis functions, but they restricted their choice of function basis to wavelet transformation, and thus can be considered as a special case of our approach. [23] discussed a method similar to what we propose, but they did not give any guidance on how to choose a function basis or offer any theoretical guarantees. These approaches can be treated as special cases within our framework with a specific choice of function basis. In addition to linear function-on-function regression, [51] and [56] studied functional additive models, but did not give theoretical results on variable selection in the high-dimensional setting.

While finishing this paper, we were made aware of concurrent work by [58] that also uses a neighborhood selection approach to estimate a graphical model

from functional data. There are several differences between our approaches. First, when estimating the neighborhood of a particular node, [58] project each observed function to its own eigenfunction basis. We consider several alternative choices for the functional basis expansion and suggest a different approach: when estimating the neighborhood of a node  $j$ , we project all observed functions onto the same basis—typically the eigenbasis of the  $j$ -th node. However, we may use different bases when estimating different neighborhoods. Second, [58] assume a slightly more general setting where the FPCA scores arise from a non-parametric additive regression model, whereas the scores arise from a linear model in our work. However, it is unclear what the class of joint distributions is with conditional means that satisfy the non-linear additive structure. The Gaussian setting considered by our work seems like the only commonly used distribution to have a conditional mean that satisfies the additive structure. Finally, to show the consistency of graph recovery, [58] require that the truncation dimension grows with  $n$ ; while our result shows non-asymptotic error bounds for any fixed truncation dimension and chooses the truncation dimension to satisfy a criterion that is independent of the sample size. Our theoretical analysis also relies on weaker assumptions and has better convergence rates. See Section 4.1 for a detailed comparison between our theoretical results and the ones of [58].

**1.2. Notation**

Let  $[n]$  denote the set  $\{1, 2, \dots, n\}$ . For a set  $S$ , we use  $|S|$  to denote its cardinality.

We use bold lower case letters (e.g.,  $\mathbf{a}$  and  $\mathbf{b}$ ) to denote vectors and bold upper case letters (e.g.,  $\mathbf{A}$  and  $\mathbf{B}$ ) to denote matrices. For a vector  $\mathbf{a} \in \mathbb{R}^n$ , let  $\|\mathbf{a}\|_q$  denote its  $l_q$ -norm,  $q \in [1, \infty)$ , with the usual extension to  $q = 0$  and  $q = \infty$ . For a set of indices  $I \subseteq [n]$ , we use  $\mathbf{a}_I$  to denote the vector in  $\mathbb{R}^n$  with  $a_{I,i} = v_i$  for all  $i \in I$  and  $a_{I,i} = 0$  for all  $i \notin I$ . Let  $\mathcal{G} = \{G_1, G_2, \dots, G_{N_{\mathcal{G}}}\}$  be a partitioning of the set  $[n]$  into a set of  $N_{\mathcal{G}}$  disjoint groups. The mixed norm  $\|\cdot\|_{1,q}$  is defined as  $\|\mathbf{a}\|_{1,q} = \sum_{t=1}^{N_{\mathcal{G}}} \|\mathbf{a}_{G_t}\|_q$ . For two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , we use  $\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^n a_i b_i$  to denote their inner product.

For a symmetric matrix  $\mathbf{B}$ , we use  $\rho_{\max}(\mathbf{B})$  to denote its largest eigenvalue,  $\rho_{\min}(\mathbf{B})$  to denote its smallest eigenvalue, and  $\text{tr}(\mathbf{B})$  to denote its trace. For a matrix  $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$ , we use  $\text{vec}(\mathbf{A})$  to denote the vector in  $\mathbb{R}^{n_1 n_2}$  formed by stacking the columns of  $\mathbf{A}$ . For two matrices  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and  $\mathbf{B} \in \mathbb{R}^{r \times s}$ ,  $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{nr \times ms}$  denotes their Kronecker product, with  $(\mathbf{A} \otimes \mathbf{B})_{ik,jl} = A_{ij} B_{kl}$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$ , we use  $\|\mathbf{A}\|_q$  to denote its operator norm,  $q \in [1, \infty]$ , that is,  $\|\mathbf{A}\|_q = \sup_{\mathbf{v} \in \mathbb{R}^{n_2}: \|\mathbf{v}\|_q=1} \|\mathbf{A}\mathbf{v}\|_q$ . Thus,  $\|\mathbf{A}\|_2$  denotes the maximum singular value of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n_2} \sum_{i=1}^{n_1} |A_{ij}|$ , and  $\|\mathbf{A}\|_{\infty} = \max_{1 \leq i \leq n_1} \sum_{j=1}^{n_2} |A_{ij}|$ . We use  $\|\mathbf{A}\|_F$  to denote the Frobenius norm of  $\mathbf{A}$ , that is,  $\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^T \mathbf{A})$ . We use  $\|\mathbf{A}\|_{\infty}$  to denote the elementwise maximum absolute value of  $\mathbf{A}$ , that is,  $\|\mathbf{A}\|_{\infty} = \max_{i,j} |A_{ij}|$ .

For a real-valued differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we use  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$  to denote its gradient at a point  $\mathbf{x} \in \mathbb{R}^n$ . For a closed interval  $\mathcal{T} \subseteq \mathbb{R}$ , we

define  $\mathcal{L}^2(\mathcal{T})$  to be the Hilbert space of square-integrable real-valued functions defined on domain  $\mathcal{T}$ , where for  $f, g \in \mathcal{L}^2(\mathcal{T})$ , we use  $\langle f, g \rangle = \int_{\mathcal{T}} f(t)g(t)dt$  to denote their inner-product and  $\|f\| = (\int_{\mathcal{T}} f^2(t)dt)^{1/2}$  to denote the  $L_2$ -norm of  $f$ . For a bivariate function  $B(t', t)$  defined on  $\mathcal{T} \times \mathcal{T}$ , we use  $\|B\|_{\text{HS}} = \|B(t', t)\|_{\text{HS}} = (\int_{\mathcal{T} \times \mathcal{T}} B^2(t', t)dt'dt)^{1/2}$  to denote its Hilbert-Schmidt norm. We use  $\mathbf{f}(\cdot) = (f_1(\cdot), f_2(\cdot), \dots, f_n(\cdot))^{\top}$  to denote a vector with function entries.

For any two Hilbert spaces  $\mathbb{H}$  and  $\mathbb{G}$ , we define  $\mathcal{B}(\mathbb{H}, \mathbb{G})$  as the set of bounded linear operators and  $\mathcal{B}_{\text{HS}}(\mathbb{H}, \mathbb{G})$  as the set of Hilbert-Schmidt operators from  $\mathbb{H}$  to  $\mathbb{G}$ . Thus,  $\mathcal{B}_{\text{HS}}(\mathbb{H}, \mathbb{G}) \subseteq \mathcal{B}(\mathbb{H}, \mathbb{G})$ . For  $h \in \mathbb{H}$  and  $g \in \mathbb{G}$ , the outer product  $g \otimes h : \mathbb{H} \mapsto \mathbb{G}$  is the rank-one linear operator  $(g \otimes h)(h') := \langle h, h' \rangle_{\mathbb{H}} g$ . When  $\mathbb{H} = \mathbb{G}$ , we let  $\mathcal{B}(\mathbb{H}) = \mathcal{B}(\mathbb{H}, \mathbb{H})$  and  $\mathcal{B}_{\text{HS}}(\mathbb{H}) = \mathcal{B}_{\text{HS}}(\mathbb{H}, \mathbb{H})$ . For any operator  $\mathcal{T} : \mathbb{H} \mapsto \mathbb{G}$ , we use  $\text{ran}(\mathcal{T}) = \{\mathcal{T}(h) : h \in \mathbb{H}\} \subseteq \mathbb{G}$  to denote its range. We denote the adjoint operator of  $\mathcal{T}$  [22, Definition 3.3.2] by  $\mathcal{T}^*$  and the Moore–Penrose inverse or pseudo inverse [22, Definition 3.5.7] of  $\mathcal{T}$  by  $\mathcal{T}^{\dagger}$ . We say that an orthonormal sequence  $\{e_n\}_{n \geq 1}$  in a Hilbert space  $\mathbb{H}$  is called an orthonormal basis or a complete orthonormal system (CONS) if  $\overline{\text{span}\{e_n\}} = \mathbb{H}$ .

For any finite number of Hilbert spaces  $\mathbb{H}_1$ – $\mathbb{H}_n$ , we define their Cartesian product space  $\mathbb{H}_1 \oplus \dots \oplus \mathbb{H}_n \triangleq \oplus_{i=1}^n \mathbb{H}_i$  as  $\{(f_1, \dots, f_n) : f_i \in \mathbb{H}_i\}$ , with endowed inner product defined by  $\langle \mathbf{f}, \mathbf{g} \rangle = \sum_{i=1}^n \langle f_i, g_i \rangle$  for all  $\mathbf{f} = (f_1, \dots, f_n), \mathbf{g} = (g_1, \dots, g_n) \in \oplus_{i=1}^n \mathbb{H}_i$ . We then have  $\oplus_{i=1}^n \mathbb{H}_i$  to also be a Hilbert space. For  $\mathbf{f} = (f_1, \dots, f_n) \in \oplus_{i=1}^n \mathbb{H}_i$ , we use  $\mathbf{f}_{-j}$  to denote  $(f_1, \dots, f_{j-1}, f_{j+1}, \dots, f_n) \in \oplus_{i=1, i \neq j}^n \mathbb{H}_i$  for any  $j \in [n]$ . When  $\mathbb{H}_i = \mathbb{H}$  for all  $i \in [n]$ , we denote  $\oplus_{i=1}^n \mathbb{H}_i$  by  $\mathbb{H}^n$ .

For any two sequences  $\{a_n\}_{n \geq 1}$  and  $\{b_n\}_{n \geq 1}$ , we use  $a_n \lesssim b_n$  or  $a_n = O(b_n)$  ( $a_n \gtrsim b_n$  or  $a_n = \Omega(b_n)$ ) to denote that there exists a constant  $c \geq 0$  such that  $a_n \leq c \cdot b_n$  ( $a_n \geq c \cdot b_n$ ) for  $n$  large enough. Similarly, we use  $a_n = \tilde{O}(b_n)$  to ignore any log terms asymptotically, that is,  $a_n = \tilde{O}(b_n)$  if  $a_n = O(b_n \log^k b_n)$  for some  $k \geq 0$ . In this paper, we use  $\tilde{O}(\cdot)$  to ignore log terms of sample size, but we keep log terms of other quantities such as the number of vertices and the dimension of a truncated function.

### 1.3. Outline of the paper

The rest of the paper is organized as follows. In Section 2, we introduce the functional graphical model and our methodology to estimate the graph structure. In Section 3, we discuss the optimization algorithm used to compute the estimator. We develop theoretical guarantees for our approach in Section 4. Results on simulated and real data are reported in Section 5 and Section 6, respectively<sup>1</sup>. We conclude the paper with a discussion in Section 7.

<sup>1</sup>Code and data to replicate the results in this paper is available in the Supplementary Material of the paper [70] or at: [https://github.com/PercyZhai/FGM\\_Neighborhood](https://github.com/PercyZhai/FGM_Neighborhood).

## 2. Methodology

In this section, we briefly review the functional graphical model in Section 2.1. We introduce a neighborhood selection procedure for estimating the functional graphical in Section 2.2 and discuss a practical implementation in subsequent subsections.

### 2.1. Functional graphical model

Let  $\mathcal{T} \subseteq \mathbb{R}$  be a closed interval, and  $\mathbb{H} \subseteq \mathcal{L}^2(\mathcal{T})$  be a Hilbert subspace of  $\mathcal{L}^2(\mathcal{T})$ . Since  $\mathcal{L}^2(\mathcal{T})$  is a separable Hilbert space,  $\mathbb{H}$  and  $\mathbb{H}^n$  are also separable Hilbert spaces for any  $1 \leq n < \infty$ . Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and  $\mathbf{g} : \Omega \mapsto \mathbb{H}^p$  be a Gaussian random element; that is, for any  $\mathbf{h} \in \mathbb{H}$ , we have  $\langle \mathbf{g}, \mathbf{h} \rangle$  be a real-valued Gaussian random variable<sup>2</sup>. We can express  $\mathbf{g}$  as  $\mathbf{g}(\omega, t) = (g_1(\omega, t), \dots, g_p(\omega, t))$ , where  $(\omega, t) \in \Omega \times \mathcal{T}$ , and for all  $\omega \in \Omega$  and  $j \in [p]$ , we have  $g_j(\omega, \cdot) \in \mathbb{H}$  to be a function with domain  $\mathcal{T}$ . In the rest of the paper, we suppress the dependency on  $\omega$ , and denote  $\mathbf{g}(t) = (g_1(t), \dots, g_p(t))$  for  $t \in \mathcal{T}$ , or sometimes even suppress the dependency on  $t$  and let  $\mathbf{g} = (g_1, \dots, g_p)$ .

To simplify the discussion, we assume that  $\mathbf{g}$  is zero mean, that is,  $\mathbb{E}[g_j] = 0$  for all  $j \in [p]$ . Furthermore, based on the Gaussian property, we have  $\mathbb{E}[\|\mathbf{g}\|^2] < \infty$ . Thus, for all  $j \in [p]$ , we can define the covariance operator of  $g_j$  as

$$\mathcal{K}_j := \mathbb{E}[g_j \otimes g_j], \tag{1}$$

and we have  $\mathcal{K}_j \in \mathcal{B}_{\text{HS}}(\mathbb{H})$ . In addition, for any index set  $I, I_1, I_2 \subseteq [n]$ , we define

$$\mathcal{K}_I := \mathbb{E}[(g_j)_{j \in I} \otimes (g_j)_{j \in I}], \quad \mathcal{K}_{I_1, I_2} := \mathbb{E}[(g_j)_{j \in I_1} \otimes (g_j)_{j \in I_2}].$$

Furthermore, following [52], we define the conditional cross-covariance function as

$$C_{jl}(t', t) = \text{Cov}(g_j(t'), g_l(t) \mid g_k(\cdot), k \neq j, l).$$

Let  $G = (V, E)$  denote an undirected graph where  $V = [p]$  is the set of vertices or nodes, and  $E \subset V^2$  is the set of edges. The edge set  $E$  encodes the pairwise Markov property of  $\mathbf{g}$  [32] if

$$E = \{(j, l) \in V^2 : j \neq l \text{ and } g_j \not\perp g_l \mid \{g_k\}_{k \neq j, l}\}. \tag{2}$$

Let  $\mathbf{g}_i(\cdot) = (g_{i1}(\cdot), \dots, g_{ip}(\cdot))$  be a random copy of  $\mathbf{g}(\cdot)$ . The goal of this work is to estimate the set of edges  $E$  when given  $n$  i.i.d. random copies  $\{\mathbf{g}_i(\cdot)\}_{i=1}^n$ . [52] proposed to estimate  $E$  using a functional graphical lasso procedure. In contrast, we propose a neighborhood selection approach detailed in the following section. In the following, we define the neighborhood of node  $j$  as

$$\mathcal{N}_j := \{k : (j, k) \in E\}. \tag{3}$$

---

<sup>2</sup>The existence and construction of Gaussian random elements taking values in any separable Hilbert space is shown as following: By Example 1.25 of [24], we can construct Gaussian random elements taking values in  $l^2$  space, that is, the space of square summable sequences. The desired conclusion then follows from the fact that any separable Hilbert space is isometrically isomorphic to  $l^2$  [22, Theorem 2.4.17].



## 2.2. Functional neighborhood selection

We develop a neighborhood selection procedure to estimate the functional graphical model. The neighborhood selection approach can be traced back to [5] and was further developed for Gaussian graphical models in a high-dimensional setting by [42]. Specifically, [42] estimated the conditional independence graph for vector-valued data  $\mathbf{X} = (X_k)_{k \in [p]}$  drawn from a multivariate Gaussian. Properties of the multivariate Gaussian ensure that for each  $j \in [p]$ , there exist  $\{\beta_{jk}\}_{k \neq j}$  such that

$$X_j = \sum_{k \neq j} \beta_{jk} X_k + \varepsilon_j, \quad (4)$$

where  $\varepsilon_j$  is normally distributed and independent of all  $X_k$ ,  $k \neq j$ . By (4), it is clear that  $\mathcal{N}_j$  is equivalent to the set  $\{k \in [p] \setminus \{j\} : \beta_{jk} \neq 0\}$ . Thus, [42] use the variables selected from a penalized regression of  $X_j$  onto all other variables to estimate  $\mathcal{N}_j$ ; specifically,  $\hat{\mathcal{N}}_j = \{k \in [p] \setminus \{j\} : \hat{\beta}_{jk} \neq 0\}$ . After estimating each neighborhood, they combine the estimates into a single estimate of the entire graph  $G$ .

Our first contribution is to show that an analogous representation to (4) also holds for  $\mathbf{g}$  under mild conditions. We start by considering the conditional expectation  $\mathbb{E}[g_j | \mathbf{g}_{-j}]$  for  $j \in [p]$ . By Doob–Dynkin representation [25, Lemma 1.13], we have a measurable map  $\mathcal{B}_j : \mathbb{H}^{p-1} \mapsto \mathbb{H}$  such that  $\mathbb{E}[g_j | \mathbf{g}_{-j}] = \mathcal{B}_j(\mathbf{g}_{-j})$  almost surely. Due to the Gaussianity of  $\mathbf{g}$ , we have  $\mathcal{B}_j \in \mathcal{B}(\mathbb{H}^{p-1}, \mathbb{H})$ , and  $e_j := g_j - \mathbb{E}[g_j | \mathbf{g}_{-j}]$  to be Gaussian and independent of  $\mathbf{g}_{-j}$  [27]. For the purposes of this paper, we require  $\mathcal{B}_j$  to be in a more narrow class, namely the class of Hilbert-Schmidt operators. Therefore, we make the following assumption.

**Assumption 2.1.** For all  $j \in [p]$ , we assume that  $\mathcal{B}_j \in \mathcal{B}_{HS}(\mathbb{H}^{p-1}, \mathbb{H})$ .

The intuition of the requirement for Assumption 2.1 is associated with the infinite-dimensional nature of functional data. To characterize  $\mathcal{B}_j$  in general, one will need to estimate an infinite number of parameters, which is prohibitive with a finite sample size. For this reason, any practical solution must approximate  $\mathcal{B}_j$  with a finite-dimensional truncation. Since any linear bounded operator between two finite-dimensional Hilbert spaces is congruent to a matrix that has a bounded Hilbert-Schmidt (Frobenius) norm, Assumption 2.1 is necessary to ensure a bounded truncation error. This assumption is also made in [58]—see Assumption 4.6 therein.

To understand what kind of data generation process will satisfy Assumption 2.1, let us consider a special case. Suppose that

$$\text{ran}(\mathcal{K}_{\mathcal{N}_j, j}) \subseteq \text{ran}(\mathcal{K}_{\mathcal{N}_j}).$$

Then by Theorem 4.8 in [27] and noting that  $\mathbb{E}[g_j | \mathbf{g}_{-j}] = \mathbb{E}[g_j | g_k, k \in \mathcal{N}_j]$ , we have

$$\mathcal{B}_j = \left( \mathcal{K}_{\mathcal{N}_j}^\dagger \mathcal{K}_{\mathcal{N}_j, j} \right)^*.$$

Therefore, requiring that  $\mathcal{B}_j$  is Hilbert-Schmidt is equivalent to requiring that  $\mathcal{K}_{\mathcal{N}_j}^\dagger \mathcal{K}_{\mathcal{N}_j, j}$  is Hilbert-Schmidt. To illustrate when this condition holds, we assume that the left singular functions of  $\mathcal{K}_{\mathcal{N}_j, j}$  ordered by singular values coincide with the eigenfunctions of  $\mathcal{K}_{\mathcal{N}_j}$  ordered by eigenvalues. Let  $\{s_{1k}\}_{k=1}^\infty$  be the non-increasing singular values of  $\mathcal{K}_{\mathcal{N}_j, j}$  and  $\{s_{2k}\}_{k=1}^\infty$  be the non-increasing eigenvalues of  $\mathcal{K}_{\mathcal{N}_j}$ . Then requiring  $\mathcal{B}_j$  to be Hilbert-Schmidt will be equivalent to requiring  $\sum_{k=1}^\infty (s_{1k}/s_{2k})^2 < \infty$ . Intuitively,  $s_{1k}$  corresponds to the covariance between  $g_j$  and its neighbors  $(g_l)_{l \in \mathcal{N}_j}$  along a direction in  $\mathbb{H}^{|\mathcal{N}_j|}$ , while  $s_{2k}$  represents the variance of  $(g_l)_{l \in \mathcal{N}_j}$  along that direction. The condition that  $\mathcal{B}_j$  is Hilbert-Schmidt basically requires that the covariance between  $g_j$  and its neighbors decreases sufficiently fast compared to the decreasing speed of the variance of its neighbors. When Assumption 2.1 is violated, then regardless of the dimension of the space used for truncation, there always exists a subspace orthogonal to it, such that the projection of  $(g_l)_{l \in \mathcal{N}_j}$  onto it has small variance, but the covariance between the projection and  $g_j$  is relatively large. As a result, the behavior of  $\mathcal{B}_j$  on this subspace cannot be ignored, and thus we cannot get a good approximation of  $\mathcal{B}_j$  by using any finite-dimensional truncation.

Based on Assumption 2.1, we have a representation similar to (4) for  $\mathbf{g}$ , which we state in the following theorem.

**Theorem 2.1.** *Assume that Assumption 2.1 holds for all  $j \in [p]$ . Then for all  $j \in [p]$ , there exists  $\{\beta_{jk}(t, t')\}_{k \neq j}$  such that*

$$g_j(t) = \sum_{k \neq j} \int_{\mathcal{T}} \beta_{jk}(t, t') g_k(t') dt' + e_j(t), \quad (5)$$

where  $e_j(\cdot) \perp g_k(\cdot)$ ,  $k \neq j$ , and  $\|\beta_{jk}(t, t')\|_{HS} < \infty$ . In addition, for any sequence  $\{\phi_m\}_{m=1}^\infty$  being a CONS of  $\mathbb{H}$ , we have

$$\beta_{jk}(t, t') = \sum_{m, m'=1}^\infty b_{jk, mm'}^* \phi_m(t) \phi_{m'}(t') \quad a.e., \quad (6)$$

where

$$b_{jk, mm'}^* = \int_{\mathcal{T} \times \mathcal{T}} \beta_{jk}(t', t) \phi_m(t) \phi_{m'}(t') dt' dt. \quad (7)$$

*Proof.* See Appendix A.1.  $\square$

Although it is straightforward to postulate that such a linear representation holds for multivariate Gaussian random functions, to the best of our knowledge, we are the first to strictly prove it. When the index is clear from the context, we will remove the subscript  $j$  from  $\beta_{jk}(t, t')$ . Given the representation in (5), it is clear that  $\mathcal{N}_j$  defined in (3) is equivalent to

$$\mathcal{N}_j = \{k \in [p] \setminus \{j\} : \|\beta_{jk}\|_{HS} > 0\}. \quad (8)$$

We can thus adapt the neighborhood selection approach to functional data and seek to construct an estimate of the graph by first estimating each neighborhood.

---

**Algorithm 1:** Functional neighborhood selection
 

---

**Input:** Observed random functions  $\{g_i(\cdot)\}_{i=1}^n$   
**for**  $j \in V$  **do**  
   Estimate the projection basis  $\phi_j$  if it is not fixed in advance  
   Use (9) to calculate projection scores for all observed functions on  $\phi_j$   
   Given projection scores, solve (13)  
   Estimate  $\mathcal{N}_j$  using (14)  
**end for**  
 Combine all neighborhoods into the estimated edge set using AND/OR rule  
**Output:** Return  $\hat{E}$

---

We denote the size of the neighborhood as  $s_j = |\mathcal{N}_j|$ . To estimate the neighborhood for  $j \in V$ , we regress  $g_j$  on  $\{g_k : k \in [p] \setminus \{j\}\}$  using a penalized functional regression approach. Despite the conceptual simplicity and high level similarity to [42], there are numerous technical challenges that need to be addressed in the functional data setting, which we discuss in Section 4.

### 2.3. Vector-on-vector regression

When the observed functions  $\{g_i(\cdot)\}_{i=1}^n$  are infinite dimensional objects, the regression problem suggested by (5) cannot be solved directly. As a practical estimation procedure, we first approximate the function-on-function regression problem with a tractable finite dimensional vector-on-vector regression problem.

Suppose we seek to estimate  $\mathcal{N}_j$  for a fixed target node  $j \in [p]$ . As a first step, we represent potentially infinite dimensional functions using a finite  $M$ -dimensional basis. Let  $\phi_j = \{\phi_{jm}\}_{m=1}^\infty$  be an orthonormal basis of  $\mathbb{H}$ ; for now, we assume that it is given, and details on selecting an appropriate basis will be discussed in Section 2.4. Using the first  $M$  basis functions, we compute the projection scores for each  $k \in [p]$  and  $m \in [M]$ :

$$a_{ikm} = \langle g_{ik}, \phi_{jm} \rangle = \int_{\mathcal{T}} g_{ik}(t) \phi_{jm}(t) dt, \quad (9)$$

and form the projection score vectors  $\mathbf{a}_{i,k,M} = (a_{ik1}, \dots, a_{ikM})^\top$ . For each observed function, the scores encode the  $L_2$  projection onto the first  $M$  elements of  $\phi_j$  and  $g_{ik}(\cdot) \approx \sum_{m=1}^M a_{ikm} \phi_{jm}(\cdot)$ .

The target node  $j$ , will typically be fixed, so for ease of presentation, we assume  $j = p$ . Furthermore, we follow the commonly used regression notation and denote the random function of the target node,  $g_{ij}(\cdot)$ , as  $g_i^Y(\cdot)$  and denote the other  $p-1$  random functions as  $(g_i^{X_1}(\cdot), \dots, g_i^{X_{p-1}}(\cdot))^\top$ . We let  $a_{im}^Y = \langle g_i^Y, \phi_{jm} \rangle$  and  $a_{im}^{X_k} = \langle g_i^{X_k}, \phi_{jm} \rangle$  denote the scores for observed functions and let  $\mathbf{a}_{i,M}^Y$  and  $\mathbf{a}_{i,M}^{X_k}$  denote the vectors of scores. At times, we will also use the notation

$$\mathbf{a}_{i,M}^X = \left( (\mathbf{a}_{i,M}^{X_1})^\top, \dots, (\mathbf{a}_{i,M}^{X_{p-1}})^\top \right)^\top \in \mathbb{R}^{(p-1)M}.$$

As shown in Appendix A.2,  $\mathbf{a}_{i,M}^Y$  can be represented as

$$\mathbf{a}_{i,M}^Y = \sum_{k=1}^{p-1} \mathbf{B}_{k,M}^* \mathbf{a}_{i,M}^{X_k} + \mathbf{w}_{i,M} + \mathbf{r}_{i,M}, \quad (10)$$

where

$$\mathbf{B}_{k,M}^* = (b_{k,mm'}^*)_{1 \leq m, m' \leq M} \in \mathbb{R}^{M \times M} \quad (11)$$

is a regression matrix parameter corresponding to  $\beta_{jk}(\cdot, *)$  defined in Theorem 2.1 and

$$b_{k,mm'}^* = \int_{\mathcal{T} \times \mathcal{T}} \beta_{jk}(t', t) \phi_m(t) \phi_{m'}(t') dt' dt \quad \text{for all } m, m' \geq 1.$$

For better illustrating the proposed method, we also compare the regression matrix  $\mathbf{B}_{k,M}^*$  with the conditional covariance operator  $\mathcal{B}_j$  in Assumption 2.1. By Assumption 2.1 and Appendix A.1, we have

$$\mathcal{B}_j = \sum_{k \neq j} \mathcal{B}_{jk}, \quad \text{and } \mathcal{B}_{jk} = \sum_{m=1}^{\infty} \sum_{m'=1}^{\infty} b_{k,mm'}^* \phi_m \otimes \phi_{m'}.$$

Compared to (11), we can see that the regression matrix  $\mathbf{B}_{k,M}^*$  can be regarded as a finite-dimensional approximation of  $\mathcal{B}_{jk}$  with respect to orthonormal basis  $\{\phi_m\}_{m=1}^{\infty}$ .

Besides, we have

$$r_{im} = \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} b_{k,mm'}^* a_{im'}^{X_k},$$

$\mathbf{r}_{i,M} = (r_{i1}, \dots, r_{iM})^\top$ ,  $w_{im} = \int_{\mathcal{T}} e_{ij}(t) \phi_m(t) dt$  and  $\mathbf{w}_{i,M} = (w_{i1}, \dots, w_{iM})^\top$ . The term  $\mathbf{w}_{i,M}$  is the noise vector corresponding to  $e_{ij}(\cdot)$  defined in Theorem 2.1, and  $\mathbf{r}_{i,M}$  is a bias term, which arises due to only using the first  $M$  basis functions. More details are provided in Section 4.

Based on (10), we may define the truncated neighborhood of node  $j$  as

$$\mathcal{N}_j^M := \{k \in [p] \setminus \{j\} : \|\mathbf{B}_{k,M}^*\|_F > 0\}. \quad (12)$$

Note that in contrast to  $\mathcal{N}_j$ ,  $\mathcal{N}_j^M$  depends on the finite-dimensional objects  $\mathbf{B}_{1,M}^*, \dots, \mathbf{B}_{p-1,M}^*$ , and thus it is estimable with a finite sample size. Since for  $j \notin \mathcal{N}_j$ , we have  $\beta_{jk} = 0$  a.e., which implies that  $\|\mathbf{B}_{k,M}^*\|_F = 0$  for all  $M \geq 1$ , thus we have  $j \notin \mathcal{N}_j^M$  for all  $M \geq 1$ . This way, it is clear that  $\mathcal{N}_j^M \subseteq \mathcal{N}_j$  for all  $M \geq 1$ . On the other hand, when we choose  $M$  large enough, such that  $\|\mathbf{B}_{k,M}^*\|_F > 0$  for all  $k \in \mathcal{N}_j$ , we then have  $\mathcal{N}_j^M = \mathcal{N}_j$ .

Given  $n$  i.i.d. samples  $\{\mathbf{g}_i(\cdot)\}_{i=1}^n$ , we estimate  $\mathbf{B}_{k,M}^*$ —and subsequently  $\mathcal{N}_j^M$  and  $\mathcal{N}_j$ —using a penalized least squares approach. Let  $a_{im}^Y$ ,  $a_{im}^{X_k}$ ,  $\mathbf{a}_{i,M}^Y$ , and

$\mathbf{a}_{i,M}^{X_k}$  denote the quantities arising from the  $i$ th observed sample. We select  $\mathbf{B}_{k,M}^*$  by minimizing the following objective:

$$\hat{\mathbf{B}}_{1,M}, \dots, \hat{\mathbf{B}}_{p-1,M} \in \arg \min_{\mathbf{B}_1, \dots, \mathbf{B}_{p-1}} \left\{ \frac{1}{2n} \sum_{i=1}^n \left\| \mathbf{a}_{i,M}^Y - \sum_{k=1}^{p-1} \mathbf{B}_k \mathbf{a}_{i,M}^{X_k} \right\|_2^2 + \lambda_n \sum_{k=1}^{p-1} \|\mathbf{B}_k\|_F \right\}, \quad (13)$$

where  $\lambda_n$  is a tuning parameter. In Section 3, we propose an efficient optimization algorithm to solve (13). The challenge in giving statistical guarantees for the estimators obtained by (13) lies in the fact that  $\mathbf{r}_{i,M}$  and  $\mathbf{a}_{i,M}^{X_k}$  in (10) are correlated, so  $\mathbf{B}_{1,M}^*, \dots, \mathbf{B}_{p-1,M}^*$  are not the coefficients of the best linear unbiased estimators for predicting  $\mathbf{a}_{i,M}^Y$  by  $\mathbf{a}_{i,M}^{X_k}$ , which is the general setting assumed in the group LASSO analysis. However, when the covariance between  $\mathbf{r}_{i,M}$  and  $\mathbf{a}_{i,M}^{X_k}$  is small enough in the sense discussed in Section 4,  $\hat{\mathbf{B}}_{1,M}, \dots, \hat{\mathbf{B}}_{p-1,M}$  may still be good estimators of  $\mathbf{B}_{1,M}^*, \dots, \mathbf{B}_{p-1,M}^*$ .

Given  $\hat{\mathbf{B}}_{1,M}, \dots, \hat{\mathbf{B}}_{p-1,M}$ , the estimated neighborhood set is then

$$\hat{\mathcal{N}}_j = \{k \in [p-1] : \|\hat{\mathbf{B}}_k\|_F > \epsilon_n\}, \quad (14)$$

where the threshold  $\epsilon_n$  is a tuning parameter. Finally, the estimated edge set  $\hat{E}$  is obtained by combining the estimated neighborhoods of each node. Following [42], the edge set  $\hat{E}$  can be computed by one of the following schemes:

- **AND:** if both  $j \in \hat{\mathcal{N}}_i$  and  $l \in \hat{\mathcal{N}}_j$  hold, then  $(j, l) \in \hat{E}$ ;
- **OR:** if either  $j \in \hat{\mathcal{N}}_i$  or  $l \in \hat{\mathcal{N}}_j$  holds, then  $(j, l) \in \hat{E}$ .

To operationalize the procedure, we discuss the choice of basis functions and the choice of tuning parameters in the following two sections.

#### 2.4. Choice of basis

A key element in the above procedure is the choice of the basis  $\phi_j$ . Throughout the paper, we assume that the basis is orthonormal; if the user specifies a non-orthonormal basis, it can first be orthonormalized with a procedure such as the Gram-Schmidt algorithm (Theorem 2.4.10 of [22]).

At a high level, there are two different approaches that can be used: the basis can be fixed in advance, or the basis can depend on the data. In the first approach, one uses a known basis, which could be selected via prior knowledge, or simply a commonly used basis for which projection scores can be efficiently computed (e.g., the Fourier, B-spline, and wavelet bases). The second approach uses a basis that is determined by unobserved population quantities and needs to be estimated before computing the projection scores. For example, functional PCA (FPCA) can be used to estimate a basis [54, Chapter 8]. In the previous section, we discussed vector-on-vector regression assuming that the basis  $\phi_j$  was known a priori, and here we discuss the case where the basis must be estimated.

For a chosen node  $j \in [p]$  and any  $i \in [n]$ , suppose that we have an estimate  $\{\hat{\phi}_{jm}\}_{m \geq 1}$  of the “true” basis  $\{\phi_{jm}\}_{m \geq 1}$ . Let  $\hat{a}_{im}^Y = \langle g_i^Y, \hat{\phi}_{jm} \rangle$ ,  $\hat{a}_{im}^{X_k} = \langle g_i^{X_k}, \hat{\phi}_{jm} \rangle$ ,  $\hat{\mathbf{a}}_{i,M}^Y = (\hat{a}_{i1}^Y, \dots, \hat{a}_{iM}^Y)^\top$ , and  $\hat{\mathbf{a}}_{i,M}^{X_k} = (\hat{a}_{i1}^{X_k}, \dots, \hat{a}_{iM}^{X_k})^\top$ . Similarly to (10), we have

$$\hat{\mathbf{a}}_{i,M}^Y = \sum_{k=1}^{p-1} \mathbf{B}_{k,M}^* \hat{\mathbf{a}}_{i,M}^{X_k} + \mathbf{w}_{i,M} + \mathbf{r}_{i,M} + \mathbf{v}_{i,M}, \tag{15}$$

where the additional term  $\mathbf{v}_{i,M}$  is defined in (45) in the appendix, which arises from using  $\hat{\phi}_j$  instead of  $\phi_j$ . When  $\hat{\phi}_j$  is close to  $\phi_j$ , the error term  $\mathbf{v}_{i,M}$  should be small. See the derivation of (15) in Appendix A.2.

Based on the relationship in (15), we estimate the graph structure as in the previous section, where  $\hat{\mathbf{B}}_1, \dots, \hat{\mathbf{B}}_{p-1}$  are estimated using (13) with  $\mathbf{a}_{i,M}^Y$  and  $\mathbf{a}_{i,M}^{X_k}$  replaced by  $\hat{\mathbf{a}}_{i,M}^Y$  and  $\hat{\mathbf{a}}_{i,M}^{X_k}$ . The subsequently estimated neighborhood sets are given by (14).

The most popular data-driven basis is the FPCA basis. Recall the linear Hilbert-Schmidt covariance operator  $\mathcal{K}_j$  defined in (1), which is the integral operator with the kernel being the covariance function of  $g_j$ , that is,

$$K_{jj}(t, t') = \text{Cov}(g_j(t), g_j(t')).$$

Then there exist eigenpairs  $\{\sigma_{jm}, \phi_{jm}(\cdot)\}_{m \in \mathbb{N}}$  of  $\mathcal{K}_j$  (Theorem 7.2.6 of [22]), where  $\{\sigma_{jm}\}_{m \in \mathbb{N}}$  are the eigenvalues and  $\{\phi_{jm}(\cdot)\}_{m \in \mathbb{N}}$  are orthonormal eigenfunctions. Since the covariance operator,  $\mathcal{K}_j$ , is symmetric and positive semidefinite, we assume that  $\sigma_{j1} \geq \sigma_{j2} \geq \dots \geq 0$  without loss of generality. According to the Karhunen-Loève theorem,  $g_{ij}$  can be represented as  $g_{ij}(\cdot) = \sum_{m=1}^\infty a_{ijm} \phi_{jm}(\cdot)$ , where  $a_{ijm} = \int_{\mathcal{T}} g_{ij}(t) \phi_{jm}(t) dt \sim N(0, \sigma_{jm})$  are the FPCA scores and  $a_{ijm}$  is independent of  $a_{ijm'}$  for  $m \neq m'$  [6, Theorem 1.5]. We will refer to  $\{\phi_{jm}(\cdot)\}_{m=1}^\infty$  as the *FPCA basis*. Since the basis is orthonormal, the function  $g_{ij}^M(\cdot) = \sum_{k=1}^M a_{ijk} \phi_{jk}(\cdot)$  is the  $\mathcal{L}^2$ -projection of  $g_{ij}(\cdot)$  onto the basis spanned by the first  $M$  FPCA functions. The main advantage of this basis is that it provides the best approximation in the  $L_2$  sense when projecting a function onto a fixed number of basis functions.

Unfortunately, the FPCA basis is typically unknown, as  $K_{jj}(t', t)$  is unknown. Therefore, we first estimate the functional covariance with the empirical version:

$$\hat{K}_{jj}(t, t') = \frac{1}{n} \sum_{i=1}^n g_{ij}(t) g_{ij}(t'). \tag{16}$$

Subsequently, an eigen-decomposition of  $\hat{K}_{jj}(t, t')$  produces the estimated eigenpairs  $\{\hat{\sigma}_{jm}, \hat{\phi}_{jm}(\cdot)\}_{m=1}^M$ , which in turn can be used to estimate FPCA scores  $\hat{a}_{ijm} = \int_{\mathcal{T}} g_{ij}(t) \hat{\phi}_{jm}(t) dt$ .

[52] and [58] also use projection scores from a dimension reduction procedure. However, there are several key differences between our approach and theirs. First, although it is the most commonly used basis, we do not restrict ourselves

to the FPCA basis, and instead consider a generic basis. This provides additional flexibility and allows us to explore the effect of the chosen basis on empirical performance. See Section 4.3 for more details. Our methodology also differs in a second, more substantial way. Both [52] and [58] project each random function on its own FPCA basis and consider the resulting projection scores for all subsequent tasks. In contrast, when estimating the neighborhood of a specific node—rather than projecting each random function onto its own subspace—we project all random functions onto the same subspace. Concisely put, the subspace to estimate  $\mathcal{N}_j$ ,  $\phi_j$ , may differ from  $\phi_k$ , the subspace used to estimate  $\mathcal{N}_k$ . However, when estimating  $\mathcal{N}_j$  we use projection scores for all functions projected on a single basis  $\phi_j$ .

Intuitively, the advantage of this approach is that we can tailor the finite dimensional representation to maximize the information relevant to selecting the neighborhood of a specific node,  $\mathcal{N}_j$ . The FPCA basis for each random function maximizes the “retained information” for that random function. Although there may be significant features of  $g_{ik}(\cdot)$  that are captured by its FPCA basis, these features may not be relevant to estimate the neighborhood of a specific node  $\mathcal{N}_j$ . Ultimately, we should care more about how  $g_{ik}(\cdot)$  behaves in the subspace spanned by  $g_{ij}$ ’s FPCA basis, which captures  $g_{i,Y}$ ’s variability, rather than the subspace spanned by its own FPCA basis. We examine a theoretical justification in Section 4.3 and also illustrate the advantages in simulations.

More concretely, using a single basis for selecting  $\mathcal{N}_j$  also avoids issues of colinearity that may arise artificially. For example, suppose that

$$g_{ik}(\cdot) = \sum_{m=1}^{\infty} a_{ikm} \phi_{km}(\cdot) \quad \text{and} \quad g_{il}(\cdot) = \sum_{m=1}^{\infty} a_{ilm} \phi_{lm}(\cdot)$$

have eigenfunctions  $\{\phi_{km}(\cdot)\}_{m \geq 1}$  and  $\{\phi_{lm}(\cdot)\}_{m \geq 1}$  that differ drastically, but  $a_{ikm}$  and  $a_{ilm}$  are highly correlated. When estimating  $\mathcal{N}_j$  using the projection scores from the FPCA basis of  $k$  and  $l$ , this would result in a poorly conditioned problem that may violate the irrepresentability condition (e.g., Assumption 4.8 in [58] or Condition 5 in [52]), despite the fact that the actual random functions  $g_{ik}(\cdot)$  and  $g_{il}(\cdot)$  are not difficult to distinguish. Projecting  $g_{ik}(\cdot)$  and  $g_{il}(\cdot)$  onto the same basis— $\phi_j$ —would avoid this concern, and the resulting projection scores would only be colinear if the actual random functions are similar and the problem is intrinsically hard.

While our methodology and theory allow for any orthonormal basis, we show both theoretically and in simulations that a well-chosen basis can improve performance. When choosing a basis, there are at least two objectives to consider. First, we want to minimize the covariance between  $\mathbf{r}_{iM}$  and  $\{\mathbf{a}_{i,M}^{\mathbf{X}_k}\}_{k \in [p-1]}$  in (10). Second, we want to maximize the signal strength  $\min_{k \in \mathcal{N}_j} \|\mathbf{B}_{k,M}^*\|_F$ . In general, simultaneously achieving these two objectives is practically infeasible. Thus, in practice, we focus on achieving at least one of the two. Achieving the first objective is generally infeasible without further restrictive assumptions (see Section 4.3). Thus, in practice, we generally focus on the second objective, which will lead us to use the FPCA basis of  $g_i^Y$ , which we recommend as a default

choice. Finally, we acknowledge that our study on the choice of function basis is far from complete. One should treat our guidance as a heuristic design, and we leave more thorough studies on this topic for further research.

### 2.5. Selection of tuning parameters

There are three tuning parameters that need to be chosen to implement our algorithm: the number of basis functions used for dimension reduction,  $M$ ; the thresholding parameter from (13),  $\epsilon_n$ ; and the group lasso penalty parameter  $\lambda_n$  in (12). We now discuss how to choose them in practice.

We first discuss how to choose the number of basis functions  $M$ . We follow the same cross-validation (CV) tuning strategy as in [52]. In practice, we have access to observations  $\{(t_{ike}, h_{ike})\}_{e=1}^{E_{ik}}$ ,  $i \in [n]$  and  $k \in [p]$ , where  $h_{ike}$  is a noisy observation of  $g_{ik}(\cdot)$  at a time point  $t_{ike} \in \mathcal{T}$ . We then divide the time interval  $\mathcal{T}$  into  $J$  equal-size folds  $\mathcal{J}_1, \dots, \mathcal{J}_J$  with  $\mathcal{T} = \cup_{l=1}^J \mathcal{J}_l$ . For  $a \in [J]$ , we treat fold  $\mathcal{J}_a$  as the validation set, and the remaining  $J - 1$  folds as the training set. For a chosen node  $j \in [p]$ , if  $\phi_j$  is known, we then fit each function  $g_{ik}(\cdot)$  with an  $M$ -dimensional  $\phi_j$  basis  $\{\phi_{j1}(\cdot), \dots, \phi_{jM}(\cdot)\}$  via least-square on the observations  $\{(t_{ike}, h_{ike})\}$  where  $t_{ike} \notin \mathcal{J}_a$  to get  $\hat{g}_{ik}(\cdot)$ ; we then calculate the squared error between  $h_{ike}$  and  $\hat{g}_{ik}(t_{ike})$  on the validation set. We repeat this procedure for  $a = 1, \dots, J$  to compute the CV error and choose  $M$  that minimizes the CV error. In the case when  $\phi_j$  is unknown, we first fit  $g_{ij}(\cdot)$  on observations  $\{(t_{ije}, h_{ije})\}$  where  $t_{ije} \notin \mathcal{J}_a$  via a  $L$ -dimensional  $B$ -spline basis [54, Chapter 5] to get  $\hat{g}_{ij}(\cdot)$ , and subsequently use  $\hat{g}_{ij}(\cdot)$  to get  $\hat{\phi}_j(\cdot)$ . Next, we fit all functions  $g_{ik}(\cdot)$  by  $\{\hat{\phi}_{j1}(\cdot), \dots, \hat{\phi}_{jM}(\cdot)\}$  via least-square on the observations  $\{(t_{ike}, h_{ike})\}$  where  $t_{ike} \notin \mathcal{J}_a$  to get  $\hat{g}_{ik}(\cdot)$ . After following the same procedure to compute CV error, we then choose  $(M, L)$  simultaneously over a grid of  $M \leq L$  values and choose the pair with the lowest error.

Next, we describe the selection process for  $\epsilon_n$  and  $\lambda_n$ . When  $\lambda_n$  is large enough, all estimated coefficients  $\hat{\mathbf{B}}_k$  will be set to zero. Specifically, by Proposition 1, there exists a threshold  $\lambda_{\max, n} > 0$  that can be calculated from the data, such that for any  $\lambda_n > \lambda_{\max, n}$ , the result  $\hat{\mathbf{B}}_k = \mathbf{0}$  for all  $k \in [p - 1]$ . Thus, we only need to consider  $\lambda_n \in (0, \lambda_{\max, n}]$ . We found empirically that traditional  $K$ -fold cross-validation performs poorly in our setting. Therefore, for each  $j \in [p]$ , we select  $\lambda_n, \epsilon_n$  pair using selective cross-validation (SCV) [57].

For each value of  $\lambda_n$ , we use the entire data set to estimate

$$\hat{\mathbf{B}}_{\lambda_n} = (\hat{\mathbf{B}}_{\lambda_n, 1}, \dots, \hat{\mathbf{B}}_{\lambda_n, p-1})$$

by solving (13). Given any threshold parameter  $\epsilon_n$ , we can obtain an index set  $\hat{\mathcal{N}}_j(\lambda_n, \epsilon_n) \subseteq [p - 1]$  that indicates the blocks in  $\hat{\mathbf{B}}_{\lambda_n}$  that are large enough in terms of Frobenius norm, that is,

$$l \in \hat{\mathcal{N}}_j(\lambda_n, \epsilon_n) \quad \text{if and only if} \quad \|\hat{\mathbf{B}}_{\lambda_n, l}\|_F > \epsilon_n. \quad (17)$$

For  $l \in \hat{\mathcal{N}}_j(\lambda_n, \epsilon_n)$ , we then re-estimate  $\hat{\mathbf{B}}_k$  by minimizing the unpenalized least squares objective using the  $k$ -th-fold training set, which we denote as  $I_k$ ,



---

**Algorithm 2:** The Selective Cross-Validation (SCV) algorithm to choose  $(\lambda_n, \epsilon_n)$ .

---

**Input:**  $A^X, A^Y, j \in [p]$ ;  
**for all**  $\lambda_n$  **do**  
  Run Group Lasso ADMM on  $(A^X, A^Y)$  and obtain  $\hat{B}_{\lambda_n}$ ;  
  **for all**  $\epsilon_n$  **do**  
    Obtain  $\hat{\mathcal{N}}_j(\lambda_n, \epsilon_n)$  by (17);  
    **for**  $k \in [K]$  **do**  
      Re-estimate  $\tilde{B}_l$  for  $l \in [p-1]$  by solving (18) with the  $k$ -th-fold training set;  
      Evaluate the estimate on the  $k$ -th-fold test data using the SCV-RSS criterion;  
    **end for**  
    Calculate the mean of the criterion across all  $K$  folds;  
  **end for**  
**end for**  
Pick the  $(\lambda_n, \epsilon_n)$  pair that minimizes the mean criterion;

---

and we set  $\tilde{B}_l = 0$  for all  $l \notin \hat{\mathcal{N}}_j(\lambda_n, \epsilon_n)$ . Specifically, we obtain  $\tilde{B}_1, \dots, \tilde{B}_{p-1}$  by solving the optimization problem below:

$$\begin{aligned} \tilde{B}_1, \dots, \tilde{B}_{p-1} \in \arg \min_{B_1, \dots, B_{p-1}} & \left\{ \sum_{i \in I_k} \left\| a_{i,M}^Y - \sum_{l=1}^{p-1} B_l a_{i,M}^{X_l} \right\|_2^2 \right\}, \\ \text{s.t. } & \tilde{B}_l = 0 \text{ for all } l \notin \hat{\mathcal{N}}_j(\lambda_n, \epsilon_n). \end{aligned} \quad (18)$$

We propose an error criterion named SCV-RSS, where RSS stands for the residual sum of squares. The criterion performs well in practice and adds the BIC penalty term to the squared norm of the empirical estimation error. Let

$$\hat{\mathbf{r}}_i := a_{i,M}^Y - \sum_{l=1}^{p-1} \tilde{B}_l a_{i,M}^{X_l},$$

and SCV-RSS on the test set  $I_{\text{test}}$  is defined as

$$\text{SCV-RSS}(\lambda_n, \epsilon_n) := \sum_{i \in I_{\text{test}}} \|\hat{\mathbf{r}}_i\|_2^2 + \log(|I_{\text{test}}|) \cdot |\hat{\mathcal{N}}_j(\lambda_n, \epsilon_n)|. \quad (19)$$

We then finally choose the  $(\lambda_n, \epsilon_n)$  pair that minimizes the mean of SCV-RSS over all  $K$  folds. The pseudo-code of the procedure is given in Algorithm 2.

### 3. Optimization algorithm

We propose an optimization method to solve (13) using the alternating direction method of multipliers (ADMM) [17, 7]. Note that (13) has a composite objective structure where the objective is composed of a convex smooth loss and a convex non-smooth regularization term. This composite objective is well studied in the convex optimization literature [see, for example, Section 5.1 9].

In this section, we provide an easy-to-use practical solution. Commonly used alternative methods to solve such a composite objective include ISTA (Iterative Shrinkage-Thresholding Algorithm) and FISTA (Fast ISTA). See Section 5.1 in [9] for more details. One advantage of ADMM is that it is easy to parallelize. Therefore, it is preferable when there are several machines available and the sample size or number of vertices is large [7, Chapter 8 and Chapter 10].

The pseudo-code of our method is given in Algorithm 3 and we provide additional details below. Let

$$\mathbf{A}^{\mathbf{Y}} = \begin{bmatrix} (\mathbf{a}_{1,M}^{\mathbf{Y}})^{\top} \\ (\mathbf{a}_{2,M}^{\mathbf{Y}})^{\top} \\ \vdots \\ (\mathbf{a}_{n,M}^{\mathbf{Y}})^{\top} \end{bmatrix} \in \mathbb{R}^{n \times M}, \quad \mathbf{A}^{\mathbf{X}_k} = \begin{bmatrix} (\mathbf{a}_{1,M}^{\mathbf{X}_k})^{\top} \\ (\mathbf{a}_{2,M}^{\mathbf{X}_k})^{\top} \\ \vdots \\ (\mathbf{a}_{n,M}^{\mathbf{X}_k})^{\top} \end{bmatrix} \in \mathbb{R}^{n \times M}.$$

Consider the concatenated matrices

$$\mathbf{A}^{\mathbf{X}} = [\mathbf{A}^{\mathbf{X}_1} \quad \mathbf{A}^{\mathbf{X}_2} \quad \dots \quad \mathbf{A}^{\mathbf{X}_{p-1}}] \in \mathbb{R}^{n \times (p-1)M},$$

and

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \\ \vdots \\ \mathbf{P}_{p-1} \end{bmatrix} \in \mathbb{R}^{(p-1)M \times M}, \quad \mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 \\ \mathbf{Q}_2 \\ \vdots \\ \mathbf{Q}_{p-1} \end{bmatrix} \in \mathbb{R}^{(p-1)M \times M}.$$

Then (13) can be reformulated as:

$$\min_{\mathbf{P}, \mathbf{Q}} \frac{1}{2n} \|\mathbf{A}^{\mathbf{Y}} - \mathbf{A}^{\mathbf{X}} \mathbf{Q}\|_{\mathbb{F}}^2 + \lambda_n \sum_{k=1}^{p-1} \|\mathbf{P}_k\|_{\mathbb{F}} \quad \text{subject to } \mathbf{P} - \mathbf{Q} = \mathbf{0},$$

which can be minimized by solving a series of optimization problems. At the  $h$ 'th iteration, for all  $k \in [p-1]$ :

$$\mathbf{P}_k^{h+1} = \arg \min_{\mathbf{P}_k} \left( \lambda_n \|\mathbf{P}_k\|_{\mathbb{F}} + \frac{\rho}{2} \|\mathbf{P}_k - \mathbf{Q}_k^h + \mathbf{U}_k^h\|_{\mathbb{F}}^2 \right), \quad k \in [p-1], \quad (20)$$

$$\mathbf{Q}^{h+1} = \arg \min_{\mathbf{Q}} \left( \frac{1}{2n} \|\mathbf{A}^{\mathbf{Y}} - \mathbf{A}^{\mathbf{X}} \mathbf{Q}\|_{\mathbb{F}}^2 + \frac{\rho}{2} \|\mathbf{Q} - \mathbf{P}^{h+1} - \mathbf{U}^h\|_{\mathbb{F}}^2 \right), \quad (21)$$

$$\mathbf{U}^{h+1} = \mathbf{U}^h + \mathbf{P}^{h+1} - \mathbf{Q}^{h+1}. \quad (22)$$

Here,  $\rho$  is the penalty parameter for the augmented Lagrangian. The solution to (20) is a group soft-thresholding update of  $\mathbf{P}$ . For each  $k \in [p-1]$ ,

$$\mathbf{P}_k^{h+1} = \left[ 1 - \frac{\lambda}{\rho \|\mathbf{Q}_k^h - \mathbf{U}_k^h\|_{\mathbb{F}}} \right]_+ (\mathbf{Q}_k^h - \mathbf{U}_k^h). \quad (23)$$

The solution to (21), i.e. the update of  $\mathbf{Q}$ , is

$$\mathbf{Q}^{h+1} = \left( \frac{1}{n} (\mathbf{A}^{\mathbf{X}})^{\top} \mathbf{A}^{\mathbf{X}} + \rho \mathbf{I}_M \right)^{-1} \left( \frac{1}{n} (\mathbf{A}^{\mathbf{X}})^{\top} \mathbf{A}^{\mathbf{Y}} + \rho \mathbf{P}^{h+1} + \rho \mathbf{U}^h \right). \quad (24)$$

---

**Algorithm 3:** ADMM for functional neighborhood selection
 

---

**Input:**  $A^X$ ,  $A^Y$ , and  $\lambda_n$ ;  
 Set initial values of  $\rho^0$ ,  $P^0$ ,  $Q^0$ , and  $U^0$ ;  
**for**  $h = 0, 1, 2, \dots$  **do**  
   Update  $P^{h+1}$  by (23);  
   Update  $Q^{h+1}$  by (24);  
   Update  $U^{h+1}$  by (22);  
   Break if primal and dual residuals meet stopping criteria;  
   Update  $\rho^{h+1}$  for next round;  
**end for**  
**Output:**  $\hat{B}_k$  for  $k \in [p-1]$ .

---

Iteratively using updates (23), (24), and (22), the matrix  $P_k^h$  will eventually converge to  $P_k^*$ ,  $k \in [p-1]$ , as  $h \rightarrow \infty$  [7]. The solution of (13) is given by  $\hat{B}_k = P_k^*$ ,  $k \in [p-1]$ . The stopping criterion for the iteration process depends on the primal residual, which indicates how well the constraints are satisfied, and the dual residual, which indicates stability of updates between two consecutive iterations [7]. In our settings,  $s_1^h = P^h - Q^h$ , and  $s_2^h = Q^h - Q^{h-1}$  are the primal and dual residuals respectively. The algorithm terminates when both residuals are below their respective tolerances:

$$\|s_1^h\|_F \leq \epsilon^{\text{pri},h} \quad \text{and} \quad \|s_2^h\|_F \leq \epsilon^{\text{dual},h},$$

where

$$\begin{aligned} \epsilon^{\text{pri},h} &= \sqrt{(p-1)M^2} \epsilon^{\text{abs}} + \epsilon^{\text{rel}} (\|P^h\|_F \vee \|Q^h\|_F), \\ \epsilon^{\text{dual},h} &= \sqrt{(p-1)M^2} \epsilon^{\text{abs}} + \epsilon^{\text{rel}} \|U^h\|_F. \end{aligned}$$

The factor  $\sqrt{(p-1)M^2}$  is because the Frobenius norms are computed on  $\mathbb{R}^{(p-1)M^2}$  matrices. In the following experiments, we use  $\epsilon^{\text{abs}} = 10^{-4}$  and  $\epsilon^{\text{rel}} = 10^{-4}$  by default.

The penalty parameter  $\rho$  of the augmented Lagrangian can be adjusted adaptively. We use Strategy S3 in Table 1 of [20] with  $\varphi = 10$ ,  $\tau^{\text{incr}} = \tau^{\text{decr}} = 2$ :

$$\rho^{h+1} = \begin{cases} \tau^{\text{incr}} \rho^h, & \|s_1^h\|_2 > \varphi \|s_2^h\|_2 \\ \rho^h / \tau^{\text{decr}}, & \|s_2^h\|_2 > \varphi \|s_1^h\|_2 \\ \rho^h, & \text{otherwise,} \end{cases}$$

with  $\rho^0 = 1$ . This guarantees that the primal and dual residuals do not vary significantly across iterations and ensures stability regardless of the initial  $P^0$  and  $U^0$ .

#### 4. Theoretical properties

We now discuss the statistical properties of the estimator proposed in Section 2. In particular, we give conditions under which the neighborhood of a single variable can be consistently recovered. Using a union bound extends the guarantees

to recovery of the entire graph. First, we discuss a procedure that uses a fixed function basis, and, subsequently, we discuss a procedure that uses an estimated function basis.

Since we first consider a single node  $j$ , we assume without loss of generality that  $j = p$ . To simplify the notation, we also drop the subscript  $j$  from  $\beta_{jk}(t', t)$ ,  $\phi_j$ ,  $\phi_{j,m}$ ,  $\dots$ , in this section. By (6), we have  $\|\beta_k(t', t)\|_{\text{HS}} = \sqrt{\sum_{m,m'=1}^{\infty} (b_{k,mm'}^*)^2}$  and  $b_{k,mm'}^* = 0$  for all  $m, m' \geq 1$  when  $\|\beta_k(t, t)\|_{\text{HS}} = 0$ . Let  $\mathbf{B}_{\mathbf{k},M}^*$  be a  $M \times M$  matrix whose  $m$ -th row is  $(b_{k,m1}^*, b_{k,m2}^*, \dots, b_{k,mM}^*)$ . The scores of the “error” projected onto the function basis are denoted as  $w_{im} = \int_{\mathcal{T}} e_i(t)\phi_m(t)dt$ ,  $m \geq 1$ , and  $\mathbf{w}_{i,M} = (w_{i1}, \dots, w_{iM})^\top$ . Let  $\mathbf{r}_{i,M} = (r_{i1}, \dots, r_{iM})^\top \in \mathbb{R}^M$  denote the “bias” arising from using the first  $M$  basis elements to represent  $\beta_k(t', t)$  where  $r_{im} = \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} b_{k,mm'}^* a_{im'}^{X_k}$ . Let

$$\begin{aligned} \beta_{k,M}(t', t) &= \sum_{m,m'=1}^M b_{k,mm'}^* \phi_m(t)\phi_{m'}(t'), \\ \beta_{k,>M}(t', t) &= \sum_{m>M \text{ or } m'>M}^{\infty} b_{k,mm'}^* \phi_m(t)\phi_{m'}(t'). \end{aligned} \tag{25}$$

Then

$$\begin{aligned} \|\beta_{k,M}(t', t)\|_{\text{HS}} &= \sqrt{\sum_{m,m'=1}^M (b_{k,mm'}^*)^2}, \\ \|\beta_{k,>M}(t', t)\|_{\text{HS}} &= \sqrt{\sum_{m>M \text{ or } m'>M}^{\infty} (b_{k,mm'}^*)^2}, \end{aligned}$$

and

$$\begin{aligned} \|\beta_k(t', t)\|_{\text{HS}} - \|\mathbf{B}_{\mathbf{k},M}^*\|_{\text{F}} &= \|\beta_{k,M}(t', t) + \beta_{k,>M}(t', t)\|_{\text{HS}} - \|\mathbf{B}_{\mathbf{k},M}^*\|_{\text{F}} \\ &\leq \|\beta_{k,M}(t', t)\|_{\text{HS}} + \|\beta_{k,>M}(s, t)\|_{\text{HS}} - \|\mathbf{B}_{\mathbf{k},M}^*\|_{\text{F}} \\ &= \|\beta_{k,>M}(t', t)\|_{\text{HS}}. \end{aligned}$$

When  $M$  is large enough, then the term  $\|\beta_{k,>M}(t', t)\|_{\text{HS}}$  is small; when  $n$  is also large enough,  $\hat{\mathbf{B}}_{\mathbf{k}}^M$  is close to  $\mathbf{B}_{\mathbf{k},M}^*$ , and  $\hat{\mathcal{N}}_j$  will be a good estimator of  $\mathcal{N}_j$ .

Both  $\mathbf{w}_{i,M}$  and  $\mathbf{r}_{i,M}$  are Gaussian vectors with mean zero, and we denote their covariance matrices as  $\Sigma^w$  and  $\Sigma^r$  respectively; in addition, we define  $\Sigma^{r,w} = \text{Cov}(\mathbf{r}_{i,M}, \mathbf{w}_{i,M})$  and  $\Sigma^{w,r} = (\Sigma^{r,w})^\top$ . To simplify the notation, we drop the explicit dependence on  $M$ . Let  $\Sigma^{\mathbf{X}_k,r} = \text{Cov}(\mathbf{a}_{i,M}^{\mathbf{X}_k}, \mathbf{r}_{i,M}) \in \mathbb{R}^{M \times M}$ ,  $\Sigma^{r,\mathbf{X}_k} = (\Sigma^{\mathbf{X}_k,r})^\top$ ,  $\Sigma^{\mathbf{X}_k,\mathbf{X}_l} = \text{Cov}(\mathbf{a}_{i,M}^{\mathbf{X}_k}, \mathbf{a}_{i,M}^{\mathbf{X}_l}) \in \mathbb{R}^{M \times M}$ , and  $\Sigma^{\mathbf{X}} = (\Sigma^{\mathbf{X}_k,\mathbf{X}_l})_{1 \leq k,l \leq p-1}$  is a matrix composed of  $M \times M$ -blocks  $\Sigma^{\mathbf{X}_k,\mathbf{X}_l}$ ,  $k, l \in [p-1]$ .

The following quantities will be used to state the results:

$$\begin{aligned}\Xi_1(M) &= \max_{k \in [p-1]} \left\{ \rho_{\max}(\Sigma^w + \Sigma^r - \Sigma^{r, \mathbf{X}_k} (\Sigma^{\mathbf{X}_k, \mathbf{X}_k})^{-1} \Sigma^{\mathbf{X}_k, r}) \right\}, \\ \Xi_2(M) &= \max_{k \in [p-1]} \rho_{\max}(\Sigma^{\mathbf{X}_k, \mathbf{X}_k}), \quad \Xi_3(M) = \max_{k \in [p-1]} \text{tr}(\Sigma^{\mathbf{X}_k, \mathbf{X}_k}), \\ \Xi_4(M) &= \text{tr} \{ \Sigma^r + \Sigma^w + \Sigma^{r, w} + \Sigma^{w, r} \}, \quad \omega(M) = \max_{k \in [p-1]} \|\Sigma^{r, \mathbf{X}_k}\|_F.\end{aligned}\quad (26)$$

Note that  $\Sigma^r - \Sigma^{r, \mathbf{X}_k} (\Sigma^{\mathbf{X}_k, \mathbf{X}_k})^{-1} \Sigma^{\mathbf{X}_k, r}$  is a conditional variance of  $\mathbf{r}_i$  given  $\mathbf{a}_{i, M}^{\mathbf{X}_k}$ , so the arguments in  $\Xi_1(M)$  are always positive semidefinite. The functions  $\Xi_1(M) - \Xi_4(M)$  are used to express an upper bound on the covariance between the projection scores  $\mathbf{a}_{i, M}^{\mathbf{X}_k}$  and the error terms  $(\mathbf{r}_{i, M} + \mathbf{w}_{i, M})$ . This upper bound then provides a lower bound for the regularization parameter  $\lambda_n$ . The function  $\omega(M)$  measures the correlation of residuals  $\mathbf{r}_{i, M}$  with  $\mathbf{a}_{i, M}^{\mathbf{X}_k}$ . A large correlation implies that the problem is more difficult to solve. Finally, let

$$K_0 = \max_{k \in [p-1], m \in M} \mathbb{E}[(a_{im}^{\mathbf{X}_k})^2] = \max_{k \in [p-1], m \in M} (\Sigma^{\mathbf{X}_k, \mathbf{X}_k})_{m, m} < \infty. \quad (27)$$

The quantity  $K_0$  is used to provide an upper bound on the estimation error for the covariance matrix of  $\mathbf{a}_{i, M}^{\mathbf{X}_k}$ . Subsequently, this is used to prove a lower bound on restricted eigenvalues, which is a crucial step in proving Theorem 4.1 and Theorem 4.3.

#### 4.1. Prior fixed function basis

Let  $\sigma_{j0} = \mathbb{E}[\|g_{ij}\|^2]$ ,  $\sigma_{\max, 0} = \max_{j \in [p]} \sigma_{j0}$ , and  $\sigma_{jr} = \mathbb{E}[\|e_{ij}\|^2]$ , where  $e_{ij}$  is defined in (5). Note that  $\sigma_{jr} \leq \sigma_{j0}$ . We introduce several assumptions before stating the main results.

**Assumption 4.1.** *There exists a constant  $C > 0$  that does not depend on  $p$  such that  $\sigma_{\max, 0} \leq C$ .*

Assumption 4.1 requires that the norm of the random functions have a finite second moment that does not grow with  $p$  and is a basic requirement for functional graphical models to be well defined. Note that  $\Xi_k(M) \leq \max_{j \in [p]} \mathbb{E}[\|g_{ij}\|^2]$  for all  $k = 1, 2, 3, 4$  and any  $M$ . Thus,  $\Xi_k(M) \leq C$  for all  $M \geq 1$  and  $p \geq 1$ .

**Assumption 4.2.** *Let  $\Sigma_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} = (\Sigma^{\mathbf{X}_k, \mathbf{X}_{k'}})_{k, k' \in \mathcal{N}_j} \in \mathbb{R}^{|\mathcal{N}_j| \times |\mathcal{N}_j|}$  be the submatrix with blocks indexed by the elements of the neighborhood set  $\mathcal{N}_j$  and define*

$$\kappa = \kappa(M) = \rho_{\min} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right). \quad (28)$$

For any  $M$ , we assume that  $\kappa > 0$ . When  $\mathcal{N}_j$  is empty, we let  $\kappa = \infty$  for all  $M$ .

Assumption 4.2 requires that the projection scores of all functions in the neighborhood of node  $j$  are linearly independent. As discussed in Section 2.4,

because we project all functions onto the same basis, projection scores would only be colinear if the functions are truly difficult to distinguish.

Let  $\tau(M)$  be the relevant signal strength:

$$\tau(M) = \min_{k \in \mathcal{N}_j} \|\mathbf{B}_k^*\|_F = \min_{k \in \mathcal{N}_j} \|\beta_{k,M}\|_{\text{HS}}, \quad (29)$$

where  $\beta_{k,M}(t', t)$  is defined in (25). For any orthonormal basis, the signal strength  $\tau(M)$  is an increasing function of  $M$ . When  $\mathcal{N}_j$  is empty, we define  $\tau(M) = \infty$  for all  $M$ . Recall that we use  $s = s_j$  to denote the size of the neighborhood,  $s = |\mathcal{N}_j|$ . As discussed in Section 2.3, when  $M$  is large enough such that  $\tau(M) > 0$ , we have  $\mathcal{N}_j^M = \mathcal{N}_j$ , where  $\mathcal{N}_j^M$  is defined in (12).

**Assumption 4.3** (Signal Strength). *We assume that*

$$\frac{\omega(M)}{\sqrt{\kappa(M)}\tau(M)}$$

*is a non-increasing function of  $M$ , and*

$$\lim_{M \rightarrow \infty} 24\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}\tau(M)} < 1.$$

Assumption 4.3 requires that the function  $\omega(M)$ —a measure of bias due to truncation—must decay quickly compared to  $\sqrt{\kappa(M)}$ , which roughly measures the conditioning of the design matrix, after dividing by  $\tau(M)$ , which measures the signal strength.

We also compare our Assumption 4.3 with Assumption 4.6 of [58]. Assumption 4.6 of [58] assumes that there exists a universal constant  $C_{\min} > 0$  such that  $\kappa(M) \geq C_{\min}$  for all  $M \geq 0$ , where  $\kappa(M)$ —defined in our Assumption 3—is the minimum eigenvalue of the covariance matrix of the function scores when using an  $M$ -dimensional basis. However, by [22, Theorem 7.2.5], the covariance operator is a compact operator; furthermore, by [22, Theorem 4.2.3], we must have  $\kappa(M) \rightarrow 0$  as  $M \rightarrow \infty$  unless  $\mathbb{H}$  has finite-rank. Thus, Assumption 4.6 of [58] is equivalent to assuming that the random functions lie in a finite-dimensional space. In contrast, in our Assumption 4, instead of assuming that  $\kappa(M)$  is uniformly bounded away from 0, we study the interplay between  $\kappa(M)$ ,  $\omega(M)$  and  $\tau(M)$ . When  $\kappa(M)$  is bounded from 0, our Assumption 4 holds; however, our assumptions also allow  $\kappa(M) \rightarrow 0$  as  $M \rightarrow \infty$ , as long as it does not decrease too quickly when compared with  $s\omega^2(M)/\tau^2(M)$ . Thus, our theory can deal with infinite-dimensional random functions. In this sense, our Assumption 4 should be considered strictly weaker than Assumption 4.6 of [58].

Let

$$\nu(M) = \tau(M) - 24\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}}. \quad (30)$$

Under Assumption 4.3,  $\lim_{M \rightarrow \infty} \nu(M) > 0$ . We denote  $M^*$  as the smallest

integer such that  $\nu(M') > 0$  for all  $M' \geq M^*$ , that is,

$$\begin{aligned} M^* &= \min \{M : \nu(M') > 0 \text{ for all } M' \geq M\} \\ &= \min \left\{ M : 24\sqrt{s} \frac{\omega(M')}{\sqrt{\kappa(M')}\tau(M')} < 1 \text{ for all } M' \geq M \right\}. \end{aligned} \tag{31}$$

Let

$$\chi(n, p, M, \delta) = \frac{6\sqrt{s}}{\sqrt{\kappa(M)}} \tilde{\lambda}(n, p, M, \delta), \tag{32}$$

where

$$\begin{aligned} \tilde{\lambda}(n, p, M, \delta) &= \tilde{O} \left( c_{n,\delta} \left\{ \frac{M\sqrt{\Xi_1(M)}}{\sqrt{n}} + \sqrt{\Xi_1(M)} \sqrt{\frac{\log(p/\delta)}{n}} \right\} \right. \\ &\quad \left. + \omega(M) \left\{ \frac{M\sqrt{\log(pM^2/\delta)}}{\sqrt{n}} + \frac{M \log(pM^2/\delta)}{n} \right\} \right), \end{aligned} \tag{33}$$

and

$$c_{n,\delta} = \sqrt{\Xi_2(M)} \cdot \left( 1 + \sqrt{\frac{2}{n} \log\left(\frac{p-1}{\delta}\right)} \right) + \sqrt{\frac{\Xi_3(M)}{n}}. \tag{34}$$

The exact form of  $\tilde{\lambda}(n, p, M, \delta)$  can be found in (54) in appendix. The function  $\tilde{\lambda}(n, p, M, \delta)$  provides a theoretical guidance on how to select the regularization parameter  $\lambda_n$  and the function  $\chi(n, p, M, \delta)$  is used to provide theoretical guidance on the choice of thresholding parameter  $\epsilon_n$ .

We are now ready to state our main result on the consistency of the neighborhood selection procedure.

**Theorem 4.1** (Neighborhood Recovery with Prior Fixed Basis). *Suppose that Assumptions 2.1-4.3 hold. Furthermore, suppose  $M \geq M^*$ ,  $\lambda_n = \tilde{\lambda}(n, p, M, \delta)$ , and  $\epsilon_n = \chi(n, p, M, \delta) + 12\sqrt{s/\kappa(M)}\omega(M)$ . Fix  $\delta \in (0, 1]$ . If the sample size  $n$  satisfies*

$$\begin{aligned} n \geq \tilde{O} \left( \max \left\{ \frac{M^4 s^2 \log(p^2 M^2/\delta)}{\kappa^2(M)}, \right. \right. \\ \left. \left. \frac{s \cdot \max \{M^2, \log(p/\delta), M^2 \omega^2(M) \log(M^2 p/\delta)\}}{\kappa(M)\nu^2(M)} \right\} \right), \end{aligned}$$

then with probability at least  $1 - \delta$ , we have

$$\sqrt{\sum_{k=1}^{p-1} \left\| \hat{\mathbf{B}}_k - \mathbf{B}_k^* \right\|_F^2} \leq \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}},$$

so that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$ .

*Proof.* See Appendix A.6. □

Note that the quantities  $\kappa(M)$ ,  $\omega(M)$ ,  $\nu(M)$ ,  $s$ , and  $M$  all implicitly depend on  $j$ . One key difference between Theorem 4.1 and a typical group lasso result [Chapter 9.6 of 63] is that the error term is correlated with the covariates—recall that the projection scores  $\mathbf{a}_{i,\mathbf{M}}^{\mathbf{X}}$  are correlated with the residual  $\mathbf{r}_{i,\mathbf{M}}$  due to the finite-dimensional approximation. The effect of the correlation is captured by the function  $\omega(M)$ . When there is no truncation bias—i.e., the random functions are finite-dimensional and  $M$  is large enough—then  $\omega(M)$  will be zero.

Using a union bound, the following corollary directly follows from Theorem 4.1 and provides guarantees for recovery of the entire graph.

**Corollary 4.2** (Graph Recovery with Prior Fixed Basis). *Suppose the conditions of Theorem 4.1 hold for all nodes  $j \in [p]$ . We use  $\kappa_j(M_j)$ ,  $\omega_j(M_j)$ ,  $\nu_j(M_j)$ ,  $s_j$ , and  $M_j$  to take the place of  $\kappa(M)$ ,  $\omega(M)$ ,  $\nu(M)$ ,  $s$ , and  $M$  in Theorem 4.1 to show their dependency on  $j$  explicitly. Let  $\hat{E}$  be the estimated edge set obtained by applying either the AND or OR rule to the estimated neighborhood of each node. If the sample size  $n$  satisfies*

$$n \geq \max_{j \in [p]} \tilde{O} \left( \max \left\{ \frac{M_j^4 s_j^2 \log(p^3 M_j^2 / \delta)}{\kappa_j^2(M_j)}, \frac{s_j \cdot \max \{M_j^2, \log(p^2 / \delta), M_j^2 \omega_j^2(M_j) \log(M_j^2 p^2 / \delta)\}}{\kappa_j(M_j) \nu_j^2(M_j)} \right\} \right),$$

then  $\mathbb{P} \{ \hat{E} = E \} \geq 1 - \delta$ .

Before moving to the next section, we compare our theorems with some existing literature.

Compared with [52], we do not assume that the functional data are finite dimensional. Instead, we study the truly infinite dimensional functional data and discuss the trade-off between bias, signal strength, and conditioning of the design matrix explicitly. When Condition 2 of [52] (which is required for the correct graph recovery therein) holds, that is, when  $g_{ij}$  is  $M(n)$ -dimensional for all  $j \in [p]$  and some positive integer  $M(n)$ , then  $\omega(M(n)) = 0$  and  $\nu(M(n)) > 0$ . However, when Assumption 4.3 holds, we do not necessarily need  $g_{ij}$  to be finite dimensional. Thus, Condition 2 of [52] is strictly stronger than Assumption 4.3.

We also compare our results to those in [58]. In addition to the difference between our Assumption 4.3 and Assumption 4.6 of [58] as discussed previously, our theoretical analysis offers an explicit characterization of a pivotal threshold in  $M$ . Specifically, we necessitate  $M$  to exceed  $M^*$ —as defined in (31)—which is contingent solely on the characteristics of the functional data, rather than  $n$ . This allows our theoretical analysis to account for finite  $M$ . In contrast, [58] necessitate a sieve-type estimator where  $M$  scales with  $n$  (Assumption 4.1 (ii)). We argue that our finding is more intuitive because increasing  $M$  seeks to decrease the approximation error. If  $M$  is not large enough, the approximation error remains too large to ensure consistent graph recovery, irrespective



of how large  $n$  may be. Therefore, accurate graph recovery becomes infeasible. Conversely, once  $M$  is large enough to render the approximation error small, consistent graph recovery is achievable irrespective of the specific  $M$  chosen, as long as  $n$  is sufficiently large. In this regard,  $M^*$  encapsulates this threshold, a concept that is absent in [58]. Furthermore, when we prescribe  $M$  to scale with  $n$  and treat  $\kappa(M)$  as a constant, our result still delivers a superior rate. To highlight this, we initially presume all other parameters remain constant and only contemplate how the sample size  $n$  relates to  $M$ . In Theorem 4.1, once  $M \geq M^*$ , the dominant term becomes the first one. When  $\kappa(M)$  is constant, we have  $n = \Omega(M^4)$ . By contrast, Assumption 4.1 (ii) of [58] stipulates that  $n = \Omega(M^{2+3\beta})$ . Noting that Assumption 4.1 (i) of [58] demands  $\beta > 1$ , it requires  $n = \Omega(M^{2+3\beta}) = \Omega(M^5)$ , rendering it inferior to our rate.

Finally, to obtain consistency of neighborhood recovery and graph recovery, we take a thresholding idea by introducing a tuning threshold  $\epsilon_n$  in (14), while both [52] and [58] rely on the irrepresentability condition [71]. Although our approach requires an additional tuning parameter, the irrepresentability condition is known to be a strong assumption. Both ideas are widely used in the literature. The hard thresholding after initial group LASSO estimation has been broadly applied in high-dimensional linear regression [43] and graphical modeling [10]. The theoretically appropriate choice of  $\epsilon_n$  depends on problem parameters, typically unknown in practice. For this reason, the hard thresholding step is primarily employed for theoretical, rather than practical, purposes. Despite recognizing the gap between practice and theory, it's crucial to note that bridging this gap is a non-trivial task and remains a long-standing challenge in high-dimensional statistics. In the simulations of Section 5, we set  $\epsilon_n = 0$ . In Appendix 5.2, we empirically demonstrate how a non-zero  $\epsilon_n$  impacts practical performance. The result shows that the benefit of a nonzero  $\epsilon_n$  is not substantial, which justifies our choice in simulations. Another way to choose  $\epsilon_n$  is by cross-validation (CV) as we described in Section 2.5. See Section 5.3 for the empirical results of CV.

#### 4.2. Data-dependent function basis

We now consider the setting where the function basis used for dimension reduction is not known in advance, and instead the basis used is an estimate of some population basis. We will assume we have access to estimates satisfying the following property.

**Assumption 4.4.** *There exist constants  $c_1, c_2 > 0$  such that for all  $0 < \delta \leq 1$ , we have*

$$\mathbb{P} \left\{ \|\hat{\phi}_{jm} - \phi_{jm}\| \leq d_{jm} \sqrt{\frac{(1/c_1) \log(c_2/\delta)}{n}} \text{ for all } m \geq 1 \right\} \geq 1 - \delta, \quad (35)$$

where  $\{d_{jm}\}$ ,  $j \in [p]$ ,  $m \geq 1$ , are constants that depend on  $j, m$  and satisfy  $d_{j0}^2 = \sum_{m=1}^{\infty} d_{jm}^2 < \infty$ ,  $j \in [p]$ .

Assumption 4.4 holds when  $\phi_{jm}(t)$ 's are the FPCA eigenfunctions of  $K_{jj}(t', t)$  and  $\hat{\phi}_{jm}(t)$ 's are the estimated FPCA eigenfunctions of  $\hat{K}_{jj}(t', t)$ —see Lemma 6 and Lemma 8 in the Supplementary Material of [52]. When Assumption 4.4 holds, let  $d_{j,\max} = \max_{m \geq 1} d_{jm}$ , and

$$d_{js}(M) = \sqrt{\sum_{m=1}^M d_{jm}^2} \quad \text{for all } M \geq 1, \tag{36}$$

so that  $d_{js}(M) \leq d_{j0}$  for all  $M \geq 1$ . In addition, let

$$\Phi(M) = \sqrt{\sum_{k=1}^{p-1} \sum_{m=1}^M \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2} \tag{37}$$

and

$$\check{\chi}(n, p, M, \delta) = \frac{6\sqrt{s}}{\sqrt{\kappa(M)}} \check{\lambda}(n, p, M, \delta), \tag{38}$$

where  $\check{\lambda}(n, p, M, \delta)$  is given in (69) in appendix.

**Theorem 4.3** (Neighborhood Recovery with Data-Dependent Basis). *Suppose Assumptions 2.1-4.4, hold. Furthermore, suppose  $M \geq M^*$ ,  $\lambda_n = \check{\lambda}(n, p, M, \delta)$ , and  $\epsilon_n = \check{\chi}(n, p, M, \delta) + 12\sqrt{s/\kappa(M)}\omega(M)$ . Fix  $\delta \in (0, 1]$ . If the sample size  $n$  satisfies*

$$n \geq \tilde{O} \left( \max \left\{ \frac{M^4 s^2 \log(p^2 M^2 / \delta)}{\kappa^2(M)}, \frac{\max \{sM^2, s \log(p/\delta), sM^2 \omega^2(M) \log(M^2 p / \delta)\}}{\kappa(M) \nu^2(M)}, \frac{\max \{s^3 M^2 (\log(1/\delta))^2, s(d_{j0}^2 - d_{js}^2(M) \Phi^2(M))\}}{\kappa(M) \nu^2(M)} \right\} \right),$$

then with probability at least  $1 - \delta$ , we have

$$\sqrt{\sum_{k=1}^{p-1} \left\| \hat{\mathbf{B}}_k - \mathbf{B}_k^* \right\|_F^2} \leq \check{\chi}(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}},$$

so that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$ .

*Proof.* See Appendix A.7. □

Comparing Theorem 4.3 with Theorem 4.1, one key difference is that the regularization parameter is increased by a term that corresponds to the estimation error of the basis functions. As a result, the sample complexity also increases due to this additional error source.

Similar to before, the following corollary provides guarantees for recovery of the whole graph and directly follows from Theorem 4.3 when applying a union bound.

**Corollary 4.4** (Graph Recovery with Data-Dependent Basis). *Suppose the conditions of Theorem 4.3 hold for all nodes  $j \in [p]$ . We use  $\kappa_j(M_j)$ ,  $\omega_j(M_j)$ ,  $\nu_j(M_j)$ ,  $s_j$ , and  $M_j$  to take the place of  $\kappa(M)$ ,  $\omega(M)$ ,  $\nu(M)$ ,  $s$ , and  $M$  in Theorem 4.1 to show their dependency on  $j$  explicitly. Let  $\hat{E}$  be the estimated edge set obtained by applying either the AND or OR rule to the estimated neighborhood of each node. If the sample size  $n$  satisfies*

$$n \geq \max_{j \in [p]} \tilde{O} \left( \max \left\{ \frac{M_j^4 s_j^2 \log(p^3 M_j^2 / \delta)}{\kappa_j^2(M_j)}, \frac{\max \{s_j M_j^2, s_j \log(p^2 / \delta), s_j M_j^2 \omega_j^2(M_j) \log(M_j^2 p^2 / \delta)\}}{\kappa_j(M_j) \nu_j^2(M_j)}, \frac{\max \{s_j^3 M_j^2 (\log(1/\delta))^2, s_j (d_{j0}^2 - d_{js}^2(M_j) \Phi_j^2(M_j))\}}{\kappa_j(M_j) \nu_j^2(M_j)} \right\} \right),$$

then  $\mathbb{P} \{ \hat{E} = E \} \geq 1 - \delta$ .

Compared to [58], our theoretical analysis is more general, since we allow, but do not restrict,  $\phi_{jm}(t)$ 's to be the FPCA eigenfunctions of  $K_{jj}(t', t)$  and  $\hat{\phi}_{jm}(t)$ 's to be the estimated FPCA eigenfunctions of  $\hat{K}_{jj}(t', t)$ .

### 4.3. Theoretical guidance on the choice of function basis

We give a theoretical guide for choosing the function basis. Note that we treat the guidance in this section as a heuristic design, and leave more thorough study on this topic for further research. Our theory can successfully explain why the PSKL basis is a good choice for functional data satisfying the partial separability condition—see Section 4.3.1.

Suppose that we use  $\phi = \{\phi_m\}_{m=1}^\infty$  as a function basis to represent the data. Let

$$\Lambda(M, \phi) = \frac{\omega(M)}{\sqrt{\kappa(M)\tau(M)}}, \tag{39}$$

where  $\omega(M)$  measures the covariance between the scores of the basis elements we include and the basis elements we truncate,  $\kappa(M)$  is the minimum eigenvalue of the covariance of the scores, and  $\tau(M)$  measures the signal strength of the basis elements we include. The function  $\Lambda(M, \phi)$  appears in Assumption 4.3, and according to the previous section, a good choice of  $\phi$  should minimize  $\Lambda(M, \phi)$  for all  $M \geq 1$ . Unfortunately, minimizing  $\Lambda(M, \phi)$  is typically infeasible, as it involves unknown quantities. We motivate two approaches for selecting the function basis. First, we show that when additional assumptions hold, a basis that minimizes  $\omega(M)$  can be used. Second, we consider a more general case and suggest approximately minimizing an upper bound on  $\Lambda(M, \phi)$ .

4.3.1. Minimize  $\omega(M)$

Our first approach to choosing the function basis is to minimize  $\omega(M)$ . To achieve that, the function basis  $\phi$  should minimize the covariance between  $\{\{g_i^{X_k}, \phi_m\}_{k=1}^{p-1}\}_{m=1}^M$  and  $\{\{g_i^{X_k}, \phi_m\}_{k=1}^{p-1}\}_{m>M}$ . Although minimizing this covariance is intractable in general, under the assumption of partial separability [67], we can solve the minimization problem exactly.

**Definition 4.5** (Partial Separability). An orthonormal function basis  $\{\phi_m\}_{m=1}^\infty$  is called the partial separability Karhunen-Loève expansion (PSKL) basis if the random vectors

$$(\langle g_i^Y, \phi_m \rangle, \langle g_i^{X_1}, \phi_m \rangle, \dots, \langle g_i^{X_{p-1}}, \phi_m \rangle), \quad m \in \mathbb{N}$$

are mutually uncorrelated.

When the PSKL basis exists, then

$$\{\{g_i^{X_k}, \phi_m\}_{k=1}^{p-1}\}_{m=1}^M \quad \text{and} \quad \{\{g_i^{X_k}, \phi_m\}_{k=1}^{p-1}\}_{m>M}$$

are uncorrelated for any  $M \geq 1$ , and  $\omega(M) = 0$  for all  $M \geq 1$ . Note that  $\Lambda(M, \phi)$  is nonnegative, thus, the PSKL basis minimizes (39) when the data generating process is partially separable. [40] proposed a test to verify the partial separability assumption. When the partial separability condition holds, by Theorem 2 of [67], one can use the eigenfunctions of  $\hat{K}(t, t') = (1/p) \sum_{j=1}^p \hat{K}_{jj}(t, t')$  as estimates of the PSKL basis, where  $\hat{K}_{jj}(t, t')$  is defined in (16). However, the partial separability assumption is strong and may not hold in general settings.

4.3.2. Minimize an approximate upper bound

When the PSKL basis does not exist, we suggest choosing a function basis by approximately minimizing the following upper bound on (39). By the calculation in Appendix A.4, we have

$$\begin{aligned} & -\log \Lambda(M, \phi) \gtrsim \\ & -\log \left\{ \max_{k \in [p-1]} \|\Sigma^{r, X_k}\|_F \right\} + \frac{1}{2} \log \{ \text{tr}(\Sigma^Y) \} + \frac{1}{2} \log \left\{ \frac{\rho_{\min}(\Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X)}{\rho_{\max}(\Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X)} \right\} \\ & + \frac{1}{2} \log \left\{ \rho_{\min} \left( [\mathbf{R}^{Y, X_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-2} [\mathbf{R}^{X_k, Y}]_{k \in \mathcal{N}_j} \right) \right\}. \quad (40) \end{aligned}$$

Therefore, by maximizing the right hand side of (40), we are approximately minimizing  $\Lambda(M, \phi)$ . Unfortunately, most of the terms depend on  $\mathcal{N}_j$ , which is unknown. As a consequence, we choose to maximize  $\log \{ \text{tr}(\Sigma^Y) \}$ , which does not depend on  $\mathcal{N}_j$ . This term is maximized when the function basis  $\phi$  is the FPCA basis of  $g_i^Y$ . More intuitively, the FPCA basis of  $g_i^Y$  maximizes the signal strength of the response variable. In Section 5, we confirm by extensive simulations that the FPCA basis of  $g_i^Y$  indeed performs well by a slight margin.

## 5. Simulations

We illustrate the finite sample properties of our neighborhood selection procedure through a simulation study. We defer the wall-clock runtime analysis of different methods to Appendix C. We generate the simulated data with the following procedure. Let

$$g_{ij}(t) = (\mathbf{a}_{ij})^\top \mathbf{f}(t), \quad i \in [n], j \in [p],$$

where  $\mathbf{a}_{ij} \in \mathbb{R}^{M^*}$  and  $((\mathbf{a}_{i1})^\top, \dots, (\mathbf{a}_{ip})^\top)^\top \in \mathbb{R}^{pM^*}$  follows a mean zero Gaussian distribution with covariance matrix  $\Sigma = \Theta^{-1}$ , and  $\mathbf{f}(t)$  is a vector that contains the first  $M^*$  Fourier basis functions. We consider the following four settings for the precision matrix  $\Theta$ .

- **Model A.** (Block Banded – Full) We generate a block-banded precision matrix  $\Theta \in \mathbb{R}^{pM^* \times pM^*}$  with  $M^* = 15$ . Define a Toeplitz matrix  $\mathbf{T}$  such that  $\mathbf{T}_{jj} = 1$ , and  $\mathbf{T}_{jl} = 1/2|j - l|$  for all  $j \neq l$ . Let  $\mathbf{A} \in \mathbb{R}^{M^* \times M^*}$  be a tridiagonal matrix with  $\mathbf{A}_{kk} = 1$  and  $\mathbf{A}_{k,k+1} = \mathbf{A}_{k+1,k} = 0.5$ . All other entries of  $\mathbf{A}$  are set to 0. The blocks of the precision matrix  $\Theta$  are then given as  $\Theta_{jj} = \mathbf{T}$ ,  $\Theta_{j,j+1} = \Theta_{j+1,j} = 0.4\mathbf{A}$ , and  $\Theta_{j,j+2} = \Theta_{j+2,j} = 0.2\mathbf{I}_{M^*}$ . All remaining blocks of  $\Theta$  are set to 0.
- **Model B.** (Block Banded – Partial) We generate a partially block-banded precision matrix  $\Theta \in \mathbb{R}^{pM^* \times pM^*}$  with  $M^* = 15$ . In this setting, every alternating block of 10 nodes have similar connection pattern as in Model A, and the remaining nodes are fully isolated. Precisely,  $\Theta$  is a block diagonal matrix, with each of its  $10M^* \times 10M^*$  blocks denoted by  $\Theta^{(k)}$ ,  $k = 1, \dots, \lceil \frac{p}{10} \rceil$ . For even  $k$ , we set  $\Theta^{(k)} = \mathbf{I}_{10M^*}$ . For odd  $k$ , we set  $\Theta^{(k)} \in \mathbb{R}^{10M^* \times 10M^*}$  to be the block-banded precision matrix such that  $(\Theta^{(k)})_{jj} = \mathbf{A}$ ,  $(\Theta^{(k)})_{j,j+1} = (\Theta^{(k)})_{j+1,j} = 0.4\mathbf{A}$ , and  $(\Theta^{(k)})_{j,j+2} = (\Theta^{(k)})_{j+2,j} = 0.2\mathbf{A}$ . All remaining blocks of  $\Theta^{(k)}$  are set to 0.
- **Model C.** (Hub Model) We generate a hub-connected precision matrix  $\Theta \in \mathbb{R}^{pM^* \times pM^*}$  with  $M^* = 5$ . We generate the edge set  $E$  from a power law distribution as follows. For each node, the number of neighbors  $m$  follows a power law distribution  $p_m = m^{-\alpha}$  with  $\alpha = 2$ , and the exact neighbors are sampled uniformly. A disjoint sequence of edge sets  $E_1, \dots, E_5$  is generated from  $E$ , yielding 5 adjacency matrices  $\mathbf{G}_l, l = 1, \dots, 5$ . The detailed algorithm of this step is given in Section 8 of [67]. Then we generate  $p \times p$  precision submatrices  $\Omega_l, l = 1, \dots, 5$ , whose supports are exactly  $\mathbf{G}_l$ . First,  $p \times p$  matrices  $\tilde{\Omega}_l$  are generated:

$$(\tilde{\Omega}_l)_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } (\mathbf{G}_l)_{ij} = 0 \text{ or } i < j \\ \sim \text{Unif}(\mathcal{D}) & \text{if } (\mathbf{G}_l)_{ij} = 1 \end{cases}$$

where  $\mathcal{D} = [-\frac{2}{3}, -\frac{1}{3}] \cup [\frac{1}{3}, \frac{2}{3}]$ . Next, we rescale the rows of  $\tilde{\Omega}_l$  so that the  $\ell_2$ -norm of each row is 1. We then obtain  $\Omega_l$  by symmetrizing  $\tilde{\Omega}_l$ ; we

average  $\bar{\Omega}_l$  with its transpose, and set the diagonals to one. Let  $\Sigma_{\text{ps}} = \text{diag}(\Sigma_1, \dots, \Sigma_5)$  with  $\Sigma_l = \varpi_l \Omega_l^{-1}$ , where  $\varpi_l = 3l^{-1.8}$ . To break the partial separability condition, we define a block precision matrix  $\bar{\Omega} \in \mathbb{R}^{pM^* \times pM^*}$ , whose diagonal blocks of  $p \times p$  are  $\bar{\Omega}_{l,l} = \Omega_l$  and off-diagonal blocks are  $\bar{\Omega}_{l,l+1} = \bar{\Omega}_{l+1,l} = (\Omega_l^* + \Omega_{l+1}^*)/2$  where  $\Omega_l^* = \Omega_l - \text{diag}(\Omega_l)$ . We then calculate

$$\Sigma_{\text{nps}} = \text{diag}(\Sigma_{\text{ps}})^{\frac{1}{2}} \left( \text{diag}(\bar{\Omega})^{-\frac{1}{2}} \bar{\Omega} \text{diag}(\bar{\Omega})^{-\frac{1}{2}} \right)^{-1} \text{diag}(\Sigma_{\text{ps}})^{\frac{1}{2}}.$$

Finally, we obtain the covariance matrix as  $(\Sigma_{jl})_{st} = (\Sigma_{\text{nps},st})_{jl}$ ,  $1 \leq j, l \leq p$  and  $1 \leq s, t \leq 5$ . Finally, the precision matrix  $\Theta = (\Sigma)^{-1}$ .

- **Model D.** (Randomly Connected) This model is similar to the setting introduced in [55], but modified to fit functional data. We generate a random block sparse precision matrix  $\Theta \in \mathbb{R}^{pM^* \times pM^*}$  with  $M^* = 15$ . Each off-diagonal block  $\Theta_{jl}$ ,  $1 \leq j \neq l \leq p$  is set to  $0.5\mathbf{I}_{M^*}$  with probability 0.1, and  $\mathbf{0}_{M^*}$  otherwise. The diagonal blocks are set as  $\Theta_{jj} = \delta' \mathbf{I}_{M^*}$ , where  $\delta'$  is chosen to guarantee the positive definiteness of the precision matrix, i.e.,  $\Theta \succ 0$ . It is sufficient to choose  $\delta'$  such that it exceeds the maximum row sum of the absolute values of the off-diagonal entries of  $\Theta$ , thus the diagonal dominance ensures positive definiteness. Notice that the partial separability condition is satisfied under this model.

Models A and B are similar to Models 1 and 2 in Section 5.1 of [52], but modified so that partial separability is violated. Model C—where partial separability is also violated—is used as a counter-example in [67]. In these three models, the partial separability condition is violated, that is, it is impossible to separate the multivariate and functional aspects of the data. However, Model D satisfies the partial separability by construction.

For each setting, we fix  $n = 100$  and let  $p = 50, 100, 150$ . Each random function is observed on a grid with  $T = 100$  equally spaced points on  $[0, 1]$ . For  $T$  observed time points,  $(t_1, \dots, t_T)$ , uniformly spread on  $[0, 1]$ , the observed values are generated by

$$g_{ij}^{\text{obs}}(t_k) = \sum_{l=1}^{M^*} a_{ijl} f_l(t_k) + \epsilon_{ijk}$$

where  $\epsilon_{ijk} \sim N(0, \sigma^2)$ . In models A, B, and D, we set  $\sigma = 0.5$ , while in model C,  $\sigma^2 = 0.05 \times \sum_{l=1}^{M^*} \text{tr}(\Sigma_l)/p$ . We report the results averaged over 50 independent runs.

In all experiments, we set  $M$ , the number of principal components used to model each function, to be 5 for all nodes. This is a typical value selected by the cross-validation process described in Section 2.5. The first simulation experiment compares the performance of our proposed method to current baseline methods, including FGLasso [52] and PSKL [67]. Since the theoretically appropriate choice of the threshold  $\epsilon_n$  depends on problem parameters that are generally unknown in practice, we use  $\epsilon_n = 0$  for this part of experiment. We plot the ROC curve

as the penalty parameter  $\lambda_n$  changes. More specifically, let  $\lambda_n = \lambda_{j,n}$  explicitly denote the parameter choice for the node  $j$ . According to Proposition A.1, there exists  $\lambda_{\max}$  such that  $\hat{\mathcal{N}}_j$  is empty. Let  $\lambda_{j,n,t} = t_\lambda \cdot \lambda_{j,\max}$ , where  $t_\lambda \in [0, 1]$  is the same for all nodes. We plot the ROC curve as  $t_\lambda$  changes from 1 to 0. The second simulation experiment illustrates the performance of our method under various choices of  $\epsilon_n$ . From the comparison of ROC's under each model setting, our empirical results confirm that optimal performance is typically achieved with a non-zero  $\epsilon_n$ . The third simulation experiment is dedicated to assess the accuracy of a single selected graph. We use the SCV-RSS criterion introduced in Section 2.5 to select  $\lambda_{j,n}$  for all nodes, and then evaluate the performance by reporting the precision and recall on the specific graph that is selected.

### 5.1. Comparison with baseline methods

We compare our method with the functional Graphical lasso (FGLasso) procedure [52] and the PSKL procedure [67]. For FGLasso, we select the parameters as proposed in [52]. For PSKL, we use the package “fgm” with the default setting [67]. As Model D is partially separable, we also implemented our method using the PSKL function basis that we assume is known a priori, which we call PSKL Basis—in this case, it is Fourier basis. To demonstrate the advantages of using a single function basis when estimating  $\mathcal{N}_j$ —as explained in Section 2.4—we implemented the following two methods to estimate the FPCA scores and compared their performances. The first method, which we call FPCA- $g_X$ , projects each function onto its own FPCA basis and uses those projection scores for all subsequent tasks. The second method, which we call FPCA- $g_Y$ , projects all other functions onto the FPCA basis of  $g_j$  when selecting  $\mathcal{N}_j$ .

To compare the methods, we plot their respective ROC curves for each model and different values of  $p$ . For each value of the tuning parameters  $\lambda_n$  and  $t_\epsilon$ , we compare the estimated edge set to the true edge set. Specifically, we calculate the true positive rate  $\text{TPR}(\lambda_n, t_\epsilon) = \text{TP}(\lambda_n, t_\epsilon) / (\text{TP}(\lambda_n, t_\epsilon) + \text{FN}(\lambda_n, t_\epsilon))$  and the false positive rate  $\text{FPR}(\lambda_n, t_\epsilon) = \text{FP}(\lambda_n, t_\epsilon) / (\text{TP}(\lambda_n, t_\epsilon) + \text{TN}(\lambda_n, t_\epsilon))$ , where  $\text{TP}(\lambda_n, t_\epsilon)$ ,  $\text{FP}(\lambda_n, t_\epsilon)$ ,  $\text{TN}(\lambda_n, t_\epsilon)$ ,  $\text{FN}(\lambda_n, t_\epsilon)$  stand for the number of true positive, false positive, true negative, and false negative number of edges, respectively. Recall that we use  $t_\epsilon = 0$  for the comparison of different methods. The ROC curves are plotted by varying the penalty parameter  $\lambda_n$ , with  $\text{TPR}(\lambda_n, 0)$  on the vertical axis and  $\text{FPR}(\lambda_n, 0)$  on the horizontal axis. We also calculate the area under the ROC curve (AUC). The ROC curves are shown in Figure 1 and the average AUC is given in Table 1.

Although FPCA- $g_Y$  and FPCA- $g_X$  perform similarly, FPCA- $g_Y$  slightly outperforms FPCA- $g_X$  across all settings. Moreover, both FPCA- $g_X$  and FPCA- $g_Y$  drastically outperform FGLasso in Models A, B, and D, and slightly outperforms FGLasso in most settings under Model C. In Models A, B, and C, where partial separability does not hold, the PSKL procedure generally underperforms the other procedures. Even in Model D, where partial separability holds, PSKL has a performance that is comparable to ours when the dimension is low. We

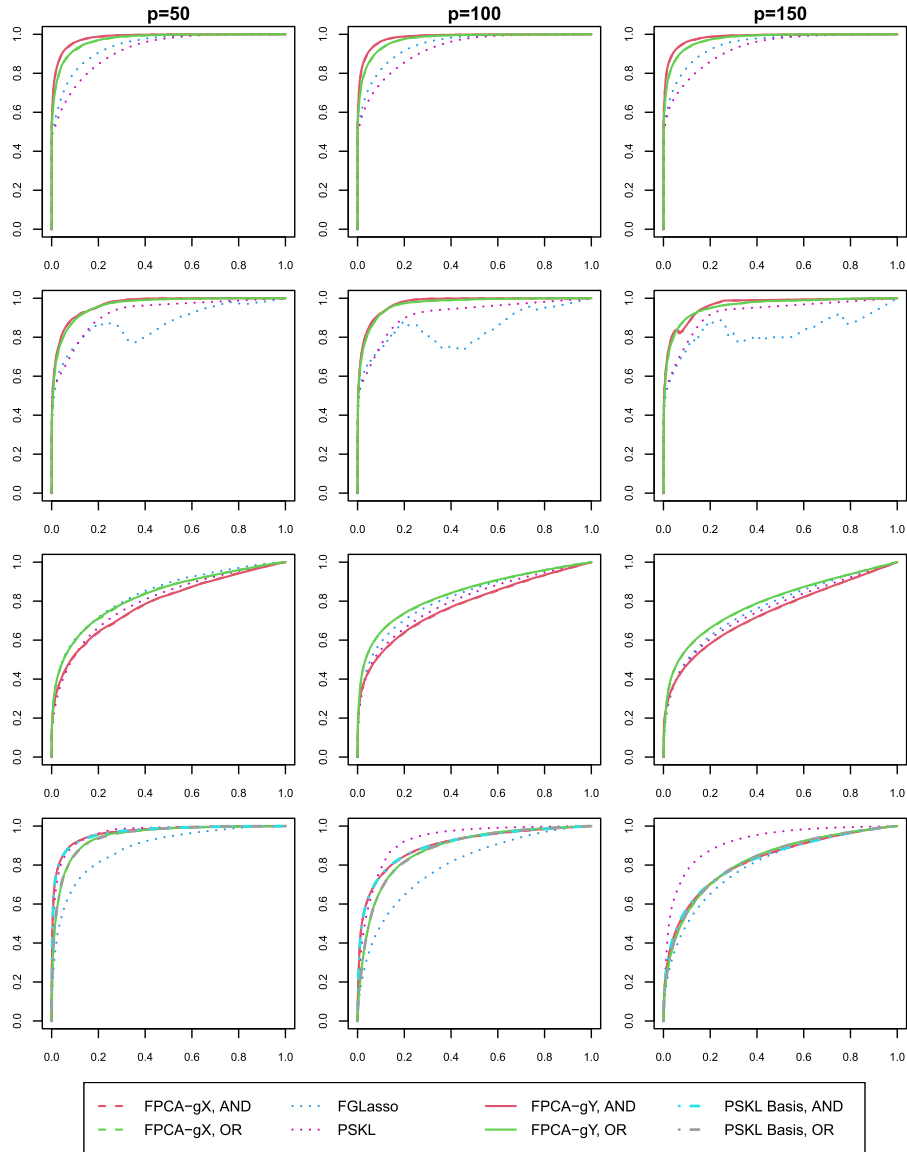


FIG 1. ROC curves for different models and values of  $p$ . From top to bottom: Models A, B, C, D. From left to right:  $p = 50, 100, 150$ . Horizontal axis: FPR; vertical axis: TPR. For FPCA-g $X$ , each function is projected onto its own FPCA basis, while FPCA-g $Y$  projects all other functions onto the FPCA basis of  $g_{ij}$ .



TABLE 1

The average AUC ROC for each method across 50 runs; the standard errors of AUC are given in the parentheses. For FPCA- $g_X$ , each function is projected onto its own FPCA basis, while FPCA- $g_Y$  projects all other functions onto the FPCA basis of  $g_{ij}$ . The maximum of each row is marked in bold.

Model	$p$	FPCA- $g_Y$ , AND	FPCA- $g_Y$ , OR	FPCA- $g_X$ , AND	FPCA- $g_X$ , OR	FGLasso	PSKL	FPCA- PSKL, AND	FPCA- PSKL, OR
A	50	<b>0.984</b> (0.004)	0.974 (0.007)	<b>0.984</b> (0.005)	0.973 (0.007)	0.942 (0.010)	0.920 (0.010)	N/A	N/A
	100	<b>0.985</b> (0.003)	0.976 (0.004)	0.984 (0.003)	0.975 (0.004)	0.947 (0.006)	0.925 (0.007)		
	150	<b>0.985</b> (0.003)	0.976 (0.003)	0.984 (0.003)	0.975 (0.003)	0.948 (0.005)	0.927 (0.007)		
B	50	<b>0.969</b> (0.008)	0.964 (0.009)	<b>0.969</b> (0.008)	0.964 (0.009)	0.806 (0.100)	0.924 (0.013)	N/A	N/A
	100	<b>0.976</b> (0.005)	0.971 (0.006)	<b>0.976</b> (0.005)	0.970 (0.006)	0.703 (0.077)	0.918 (0.021)		
	150	<b>0.965</b> (0.006)	0.961 (0.007)	0.964 (0.006)	0.960 (0.008)	0.620 (0.067)	0.924 (0.012)		
C	50	0.785 (0.035)	0.828 (0.037)	0.785 (0.035)	0.828 (0.038)	<b>0.838</b> (0.037)	0.799 (0.042)	N/A	N/A
	100	0.780 (0.040)	<b>0.839</b> (0.036)	0.777 (0.039)	0.837 (0.036)	0.822 (0.101)	0.797 (0.061)		
	150	0.740 (0.061)	<b>0.792</b> (0.053)	0.738 (0.060)	0.790 (0.053)	0.768 (0.115)	0.755 (0.077)		
D	50	<b>0.967</b> (0.012)	0.948 (0.017)	0.966 (0.013)	0.947 (0.017)	0.888 (0.081)	0.966 (0.044)	0.966 (0.013)	0.948 (0.017)
	100	0.902 (0.029)	0.882 (0.022)	0.900 (0.029)	0.881 (0.022)	0.798 (0.092)	<b>0.929</b> (0.037)	0.902 (0.030)	0.881 (0.022)
	150	0.823 (0.013)	0.824 (0.010)	0.821 (0.013)	0.822 (0.011)	0.802 (0.040)	<b>0.917</b> (0.009)	0.823 (0.013)	0.824 (0.010)

also note that in Model D, the neighborhood selection procedure that uses the “optimal” PSKL basis, which we assume is known a priori, has a very similar performance to the procedure that uses the FPCA basis and does not require prior knowledge. This suggests that using the data-selected FPCA basis is a good idea across a variety of settings.

### 5.2. The effect of $\epsilon_n$

In this section, we empirically demonstrate how  $\epsilon_n$  impacts practical performance. The experimental setting remains identical to that in Section 5. We apply our proposed method by setting  $\epsilon_n = t_\epsilon \lambda_n$ , and compute an ROC for each fixed  $t_\epsilon$  value. For each model, we select five distinct  $t_\epsilon$  values, inclusive of 0. Figure 2 provides a visually intuitive comparison and Table 2 illustrates the area under the curve (AUC) for each of the five  $t_\epsilon$  values under each setting. In most scenarios, the maximal AUC is achieved by a non-zero  $t_\epsilon$ . However, the marginal benefit of using the optimal  $t_\epsilon$  over simply setting  $t_\epsilon = 0$  is relatively insignificant in most cases. This empirical investigation justifies the practical choice of  $\epsilon_n = 0$ .

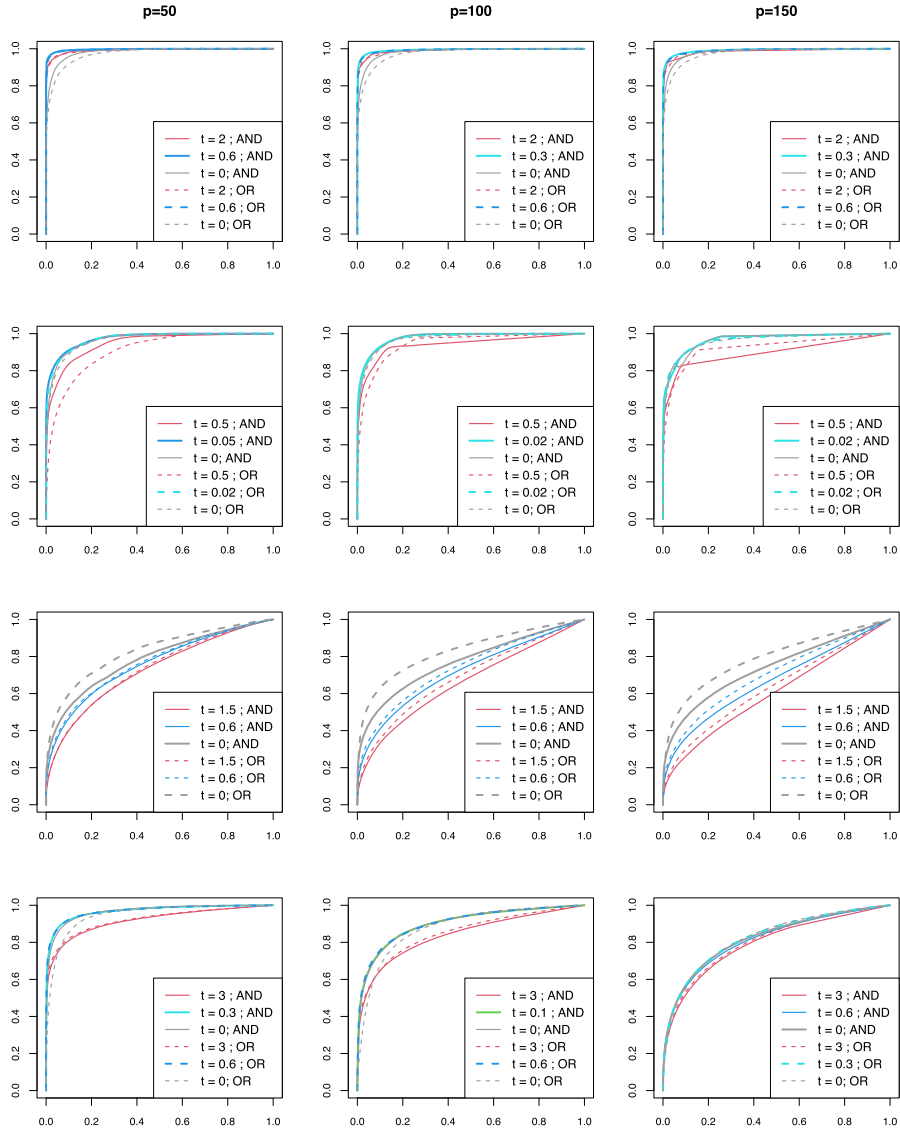


FIG 2. ROC for different models and values of  $p$  under different threshold coefficient  $t_\epsilon$  using FPCA-gx method. From top to bottom: Models A, B, C, D. From left to right:  $p = 50, 100, 150$ . Horizontal axis: FPR; vertical axis: TPR.

TABLE 2

The average AUC for each  $t_\epsilon$  across 50 runs; the standard errors of AUC are given in the parentheses. The method is FPCA-g<sub>X</sub>. The maximum of each row is marked in bold.

Model A		$t_\epsilon = 2.0$	$t_\epsilon = 0.6$	$t_\epsilon = 0.3$	$t_\epsilon = 0.1$	$t_\epsilon = 0$
p=50	AND	0.992 (0.004)	<b>0.997</b> (0.003)	0.996 (0.003)	0.993 (0.004)	0.984 (0.005)
	OR	0.991 (0.004)	<b>0.996</b> (0.003)	0.995 (0.003)	0.989 (0.004)	0.973 (0.007)
p=100	AND	0.970 (0.009)	<b>0.994</b> (0.003)	<b>0.994</b> (0.003)	0.992 (0.003)	0.984 (0.003)
	OR	0.973 (0.005)	<b>0.993</b> (0.003)	<b>0.993</b> (0.003)	0.989 (0.003)	0.975 (0.004)
p=150	AND	0.985 (0.006)	0.991 (0.003)	<b>0.992</b> (0.003)	0.990 (0.002)	0.984 (0.003)
	OR	0.987 (0.004)	<b>0.991</b> (0.002)	<b>0.991</b> (0.002)	0.987 (0.002)	0.975 (0.003)
Model B		$t_\epsilon = 0.5$	$t_\epsilon = 0.2$	$t_\epsilon = 0.05$	$t_\epsilon = 0.02$	$t_\epsilon = 0$
p=50	AND	0.901 (0.026)	0.964 (0.010)	<b>0.974</b> (0.009)	0.973 (0.009)	0.969 (0.008)
	OR	0.856 (0.020)	0.949 (0.011)	0.970 (0.009)	<b>0.971</b> (0.009)	0.964 (0.009)
p=100	AND	0.940 (0.013)	0.966 (0.009)	0.977 (0.005)	<b>0.978</b> (0.005)	0.976 (0.005)
	OR	0.943 (0.009)	0.963 (0.007)	0.975 (0.006)	<b>0.976</b> (0.006)	0.970 (0.006)
p=150	AND	0.896 (0.015)	0.946 (0.009)	<b>0.965</b> (0.007)	<b>0.965</b> (0.007)	0.964 (0.006)
	OR	0.924 (0.012)	0.952 (0.008)	0.962 (0.008)	<b>0.964</b> (0.008)	0.960 (0.008)
Model C		$t_\epsilon = 1.5$	$t_\epsilon = 0.6$	$t_\epsilon = 0.1$	$t_\epsilon = 0.05$	$t_\epsilon = 0$
p=50	AND	0.725 (0.043)	0.758 (0.041)	0.783 (0.037)	<b>0.785</b> (0.037)	<b>0.785</b> (0.035)
	OR	0.730 (0.035)	0.766 (0.036)	0.808 (0.032)	0.817 (0.033)	<b>0.828</b> (0.038)
p=100	AND	0.659 (0.079)	0.715 (0.073)	0.761 (0.060)	0.766 (0.058)	<b>0.771</b> (0.054)
	OR	0.688 (0.076)	0.738 (0.072)	0.805 (0.062)	0.818 (0.059)	<b>0.830</b> (0.061)
p=150	AND	0.602 (0.076)	0.666 (0.070)	0.723 (0.063)	0.730 (0.062)	<b>0.738</b> (0.059)
	OR	0.632 (0.081)	0.699 (0.070)	0.765 (0.067)	0.778 (0.062)	<b>0.791</b> (0.052)
Model D		$t_\epsilon = 3.0$	$t_\epsilon = 0.6$	$t_\epsilon = 0.3$	$t_\epsilon = 0.1$	$t_\epsilon = 0$
p=50	AND	0.918 (0.115)	0.966 (0.024)	<b>0.969</b> (0.015)	0.968 (0.013)	0.966 (0.013)
	OR	0.922 (0.113)	<b>0.969</b> (0.018)	0.966 (0.014)	0.956 (0.014)	0.948 (0.017)
p=100	AND	0.840 (0.104)	0.896 (0.036)	0.900 (0.031)	<b>0.901</b> (0.030)	0.900 (0.029)
	OR	0.851 (0.100)	<b>0.903</b> (0.031)	0.899 (0.027)	0.889 (0.024)	0.881 (0.022)
p=150	AND	0.794 (0.042)	0.814 (0.011)	0.818 (0.012)	0.820 (0.013)	<b>0.821</b> (0.013)
	OR	0.805 (0.034)	0.824 (0.013)	<b>0.825</b> (0.013)	0.823 (0.011)	0.822 (0.011)

TABLE 3

The average precision and recall of the graph using FPCA- $g_X$  method. The optimal  $\lambda_n$  and  $t_\epsilon$  is selected by the SCV-RSS criterion across 50 runs; the standard deviation is given in paranthesis.

Model	$p$	AND, Precision	AND, Recall	OR, Precision	OR, Recall
A	50	1.000 (0.000)	0.644 (0.050)	0.985 (0.012)	0.843 (0.037)
	100	1.000 (0.002)	0.630 (0.032)	0.970 (0.013)	0.826 (0.031)
	150	1.000 (0.001)	0.626 (0.029)	0.964 (0.010)	0.815 (0.022)
B	50	0.934 (0.071)	0.463 (0.066)	0.442 (0.120)	0.659 (0.058)
	100	0.601 (0.105)	0.528 (0.049)	0.155 (0.024)	0.746 (0.052)
	150	0.338 (0.061)	0.556 (0.038)	0.102 (0.009)	0.782 (0.032)
C	50	0.853 (0.212)	0.050 (0.033)	0.549 (0.124)	0.145 (0.048)
	100	0.902 (0.080)	0.076 (0.031)	0.646 (0.093)	0.211 (0.064)
	150	0.849 (0.085)	0.062 (0.023)	0.590 (0.070)	0.172 (0.057)
D	50	0.998 (0.014)	0.122 (0.127)	0.989 (0.030)	0.263 (0.240)
	100	0.966 (0.150)	0.034 (0.024)	0.957 (0.055)	0.114 (0.086)
	150	0.988 (0.058)	0.004 (0.009)	0.979 (0.044)	0.012 (0.031)

### 5.3. Performance of cross-validation

Practitioners may want a single graph rather than a series of graphs corresponding to different penalty and threshold parameters. Thus, we also evaluate the precision and recall of the final graph selected using the parameters obtained through selective cross-validation algorithm stated in Algorithm 2. When choosing  $\lambda_n, \epsilon_n$ , we let candidate  $\epsilon_n$ 's to be  $\epsilon_n = t_\epsilon \cdot \lambda_n$ , where  $t_\epsilon \in \{0, 0.2, 0.4, 0.8, 1.2, 1.6, 2\}$ . We denote the chosen tuning parameters as  $(\lambda_n^*, t_\epsilon^*)$ . The precision and recall of  $(\lambda_n^*, t_\epsilon^*)$  are defined as

$$\begin{aligned} \text{Precision}(\lambda_n^*, t_\epsilon^*) &= \text{TP}(\lambda_n^*, t_\epsilon^*) / (\text{TP}(\lambda_n^*, t_\epsilon^*) + \text{FP}(\lambda_n^*, t_\epsilon^*)), \\ \text{Recall}(\lambda_n^*, t_\epsilon^*) &= \text{TP}(\lambda_n^*, t_\epsilon^*) / (\text{TP}(\lambda_n^*, t_\epsilon^*) + \text{FN}(\lambda_n^*, t_\epsilon^*)). \end{aligned}$$

A larger value of precision and recall indicate better performance. The results under all models using FPCA- $g_X$  basis are shown in Table 3. From Table 3 we see that our method obtains satisfactory performance under most models, even in the high-dimensional setting. In applications where a type-I error is more costly, the AND scheme may be preferable because it enjoys a higher precision; when we want to minimize type-II errors, the OR scheme is preferred.

## 6. Data analysis

In this section, we illustrate the practical application of our method on two functional magnetic resonance imaging (fMRI) datasets. Raw brain magnetic resonance images are segmented into temporal signals for 116 regions of interest (ROIs) using the automatic anatomic labeling (AAL) parcellation approach [62]. Table 6 in Appendix D lists the names and corresponding labels of all 116 ROIs. By applying this approach, we average the signal within all ROIs to obtain 116 distinct time series, which we interpret as observations of 116 corresponding

random functions. Using the neighborhood selection procedure, we can recover the conditional independence (CI) graphs associated with different ROIs.

Recent research uncovers a hierarchical structure in brain connectivity. For instance, heteromodal areas, such as the prefrontal cortex, inferior parietal lobe, and superior temporal sulcus, project to paralimbic areas like the insula, orbitofrontal, cingulate, parahippocampal, and temporopolar regions. These, in turn, project to limbic areas, namely the amygdala and hippocampus. The latter two are the only parts of the cortex with substantial connections to the hypothalamus, a key node for homeostatic, autonomic, and endocrine aspects of the internal milieu [44]. By learning the conditional independence graph of ROIs, we gain insight into these brain connectivity patterns. Moreover, comparing conditional independence graphs from populations with and without specific neurodevelopmental conditions could yield clues about the origins of certain symptoms.

Our functional graphical models approach offers significant advantages over traditional non-functional analyses of fMRI signals. Specifically, it can detect spatio-temporal interactions among ROIs. For instance, our method can identify the influence of one node at time  $t$  on another node at a different time  $t'$ . We subsequently apply our method to two fMRI datasets: one pertaining to Autism Spectrum Disorder (ASD) and the other to Attention Deficit Hyperactivity Disorder (ADHD).

### *ASD dataset*

Autism Spectrum Disorder (ASD) is a chronic neurodevelopmental disorder associated with both sensory processing and high-level functional deficits [14]. Functional magnetic resonance imaging (fMRI) analysis provides a method for characterizing connectome anomalies in individuals with ASD.

ASD is characterized by a dissociation of a transmodal core, which combines long-distance connections from peripheral networks with primarily short-range connectivity [21]. In contrast to a neurotypical brain, which exhibits distributed functional activation patterns, an autistic brain features more regionally localized connections where selective core activation is less prominent [4].

We apply our procedure to data from the Autism Brain Imaging Data Exchange (ABIDE), a consortium that provides previously gathered fMRI data from both autism and control groups [16]. The selected samples encompass whole-brain fMRI scans from 73 ASD-diagnosed patients ( $n_{\text{autism}} = 73$ ) and 98 controls ( $n_{\text{control}} = 98$ ).<sup>3</sup> Given that  $p = 116$ , this dataset is high-dimensional. We use the time series, preprocessed by [15] using AAL parcellation, derived from the raw data.

The interpretation of the results naturally depend on the sparsity level, and the network sparsity level should be treated as a tuning parameter, which may be chosen by either domain knowledge or a data-driven approach. We initially

---

<sup>3</sup>The dataset includes fMRI measurements from eight different sites. For consistency, we only used data from New York University.

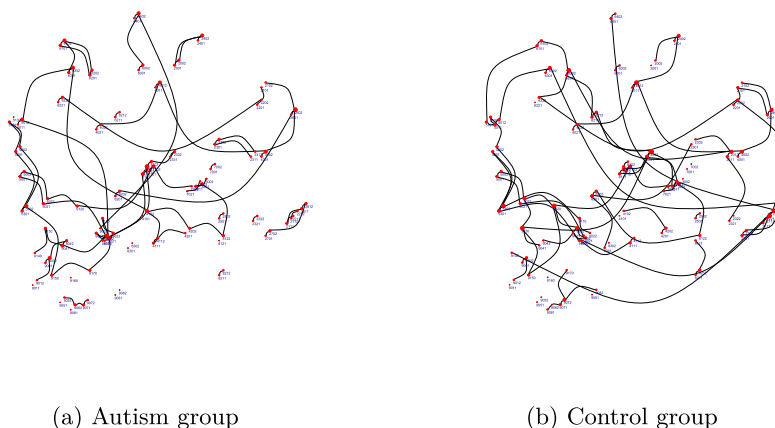


FIG 3. Brain connectome graph of autism and control groups using FPCA- $g_X$  method, obtained by SCV process

estimate the CI graphs of both autism and control groups using the SCV procedure separately, as depicted in Figure 3. Comparison of the connectivity graphs of autism and control groups reveals an overall reduction in connectivity across different brain centers in the autism group, aligning with prior findings of cortical underconnection in ASD [41]. Notably, the orbitofrontal regions (nodes 2111, 2112, 2211, 2212, 2321, 2322, 2611, 2612) appear almost isolated from the rest of the brain. This finding suggests that the orbito-frontal region, a typical paralimbic area according to [44], is less connected to limbic areas like the amygdala and hippocampus. This result is consistent with previous findings of diminished activity in the hypothalamus, leading to decreased oxytocin and vasopressin synthesis and release, which may contribute to impaired social cognition and behavior in ASD [12].

Additionally, we estimate the CI graphs of both groups under a fixed 2% sparsity following the same approach as previous analyses [52, 37] where the authors also set the network density to a small fixed level. We choose 2% because we observe that further increasing the sparsity level will induce substantially more suspicious connections in the estimated networks for both Autism and Control groups. The results are provided in Figure 4. One notable observation is the reduced rich-club connection<sup>4</sup> in the autism group. Figure 4 shows a less hierarchical brain connectome in the autism group compared to the control group. The control group exhibits more centralized connections and fewer regions without connections, while the autism group displays a more evenly distributed connection pattern across all nodes. For the control group, 22 nodes have at least 4 connections, 13 nodes have at least 5 connections, and 4 nodes have at least 6 connections. In comparison, for the autism group, 17 nodes have

<sup>4</sup>In neuroscience literature, the brain connectome structure where connections are centered around certain hub nodes is called rich-club.

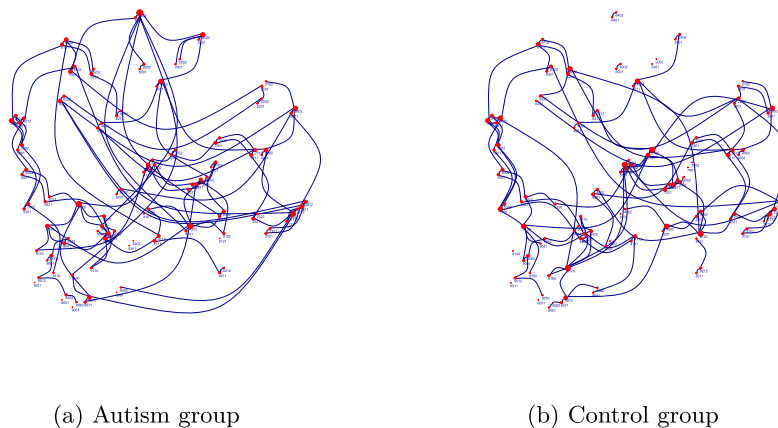


FIG 4. Brain connectome graph of autism and control groups using FPCA- $g_X$  method, with sparsity fixed to 2%

at least 4 connections, 9 nodes have at least 5 connections, and 3 nodes have at least 6 connections. Given that the total number of edges in both groups is identical, the standard deviation of the degree of all nodes is 1.52 for the control group and 1.37 for the autism group. These observations corroborate the results in [21], suggesting that ASD is associated with selective disruption in long-range connectivity, coupled with a deficit in fully activating the “rich-club.” Our findings also align with previous fMRI studies showing that individuals with ASD exhibit more spatially diffuse activations in the cerebellum’s motor-related regions [2].

Another notable observation from both Figures 3 and 4 is that the autism group displays increased connectivity in the precentral (Nodes 2001, 2002), post-central (Nodes 6001, 6002), and paracentral (Nodes 6401, 6402) regions. This observation aligns with reports by [49]. These regions are critical components of the motor control network, and abnormal activities within these areas could potentially be associated with ASD [46].

### ***ADHD dataset***

Attention Deficit Hyperactivity Disorder (ADHD) is a mental health disorder characterized by persistent issues such as difficulty maintaining attention, hyperactivity, and impulsive behavior. Functional graphical modeling may be instrumental in identifying abnormal brain connectivity associated with this condition.

We apply our procedure to data from the ADHD-200 Consortium [45]. The samples used in our analysis include whole-brain fMRI scans from 74 ADHD-diagnosed patients ( $n_{\text{ADHD}} = 74$ ) and 109 controls ( $n_{\text{control}} = 109$ )<sup>5</sup>. This

<sup>5</sup>The dataset includes fMRI measurements from eight different sites. For consistency, we

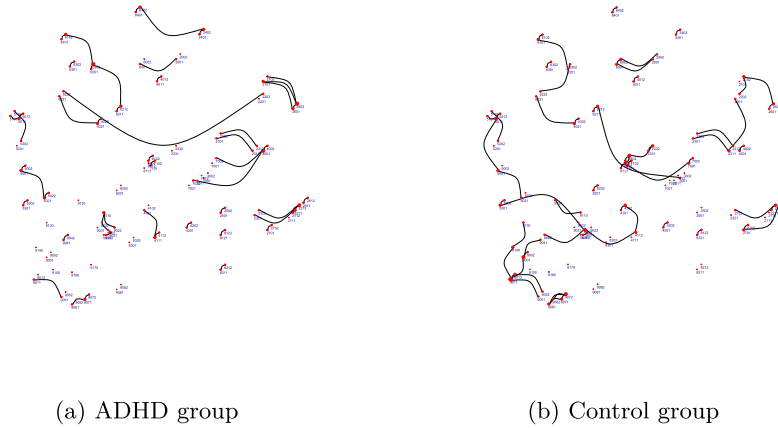


FIG 5. Brain connectome graph of ADHD and control groups using FPCA- $g_X$  method, obtained by SCV process

dataset is high-dimensional, as neither sample size exceeds  $p = 116$ . The time series preprocessed by [3] using AAL parcellation from the raw data is used in our study.

We initially estimate the CI graphs of both the ADHD and control groups using the SCV procedure separately, as demonstrated in Figure 5. We observe significantly reduced brain connectivity in the ADHD group across the entire brain network. The connectivity graph of the ADHD group in Figure 5a contains 51 edges, while the control group in Figure 5b has 62 edges. This observation aligns with the findings in [65] suggesting decreased homotopic, intrahemispheric, and heterotopic functional connectivity (i.e., disconnection) within the ADHD group. Specifically, a weaker connection is apparent within the cerebellum regions (nodes on the bottom left of Figures 5a and 5b with labels beginning with “90”) in the ADHD group. This observation is consistent with the conclusion in [11] stating that individuals with ADHD exhibit altered connectivity in cerebellum circuits, which are linked to timing disorders.

Furthermore, we estimate the CI graphs of both groups under a fixed 2% sparsity. Similar to the analysis of ASD dataset, we choose 2% because we find that further increasing the sparsity level will induce substantially more suspicious connections in the estimated network for both ADHD and Control groups. The results are shown in Figure 6. The connectivity graph of the control group in Figure 6b features several highly centralized areas, for instance, the paracentral lobule (Nodes 6401 and 6402 at the top) and prefrontal regions (e.g., Nodes 2111, 2112, 2211, 2212, 2321, 2322, 2611, 2612 on the right). The connectivity to these rich-club nodes is markedly reduced in the ADHD group. These rich-club connections are theorized to play a central role in integrating information among different brain subsystems. ADHD may be characterized by diminished struc-

---

solely utilized data from Peking University that passed quality tests.



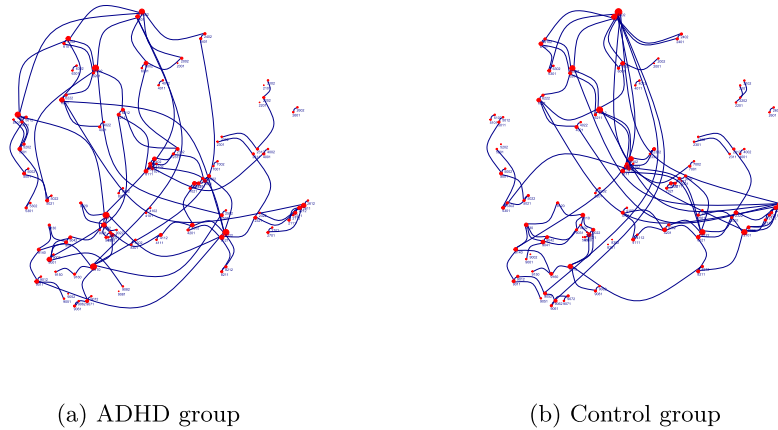


FIG 6. Brain connectome graph of ADHD and control groups using FPCA- $g_X$  method, with sparsity fixed to 2%

tural integrity of the rich-club backbone, potentially leading to a decrease in globally efficient communication capacity and altered functional brain dynamics [64]. Specifically, the diminished prefrontal activities in ADHD have been pinpointed by neuroscientific studies [30, 11]. Deficits in these regions have been associated with impairments in cognitive functions and the capacity to adapt behavior to changing circumstances flexibly [8]. Differences between the ADHD and control group are also identified in non-rich-club regions. The connections stemming from the precuneus regions (Nodes 6301 and 6302 at the top left of the graph) are markedly reduced in the ADHD group—their connections to the inferior parietal (Nodes 6101, 6102) and paracentral regions are no longer detected. The precuneus is linked to functional disturbances in regulatory control, attention, and aspects of executive function. Our observation aligns with the findings of [48], which underscore connectivity abnormalities in the precuneus among ADHD patients.

We have also applied the FGLasso method by [52] to the ADHD dataset, adjusting the connection sparsity to 2% by tuning the penalty parameter. The resulting connectome graph can be seen in Figure 7. When compared with Figure 6b, it's noticeable that FGLasso tends to generate more rich-club results. This observation aligns with simulation results wherein FGLasso exhibits relatively good performance when the underlying model features a rich-club connection structure (Model C). However, even when the underlying model has minimal rich-club structure (e.g., Models A, B, and D), FGLasso still tends to impose a rich-club structure, leading to subpar performance. As a consequence, while the FGLasso method is effective in identifying the most active regions in the connectome, it may result in a biased conclusion if such a rich-club structure is not present in reality.

For instance, according to Figure 7, the visual cortex region of the ADHD patients (Nodes 5001, 5002, 5011, 5012, 5021, 5022, 5101, 5102, 5201, 5202, 5301,

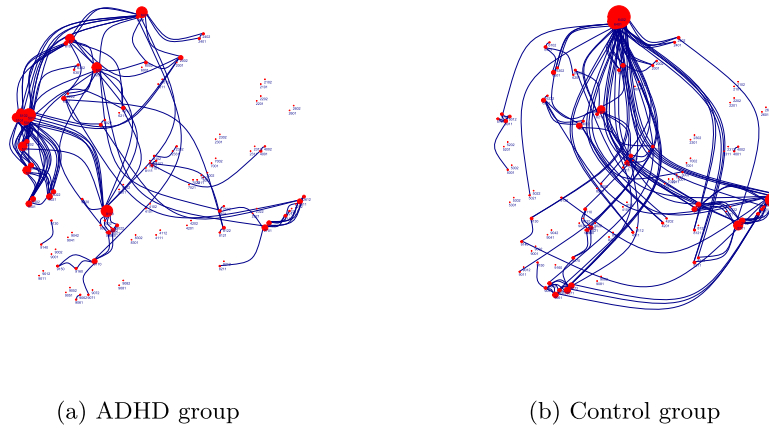


FIG 7. Brain connectome graph of ADHD and control groups using FGLasso method [52], with sparsity fixed to 2%

5302) appears to be densely connected to other brain regions. However, considering that the fMRI dataset we use is gathered during a resting state and that the visual cortex is primarily dedicated to visual functions [19], such a connection pattern within the ADHD group is unexpected. Furthermore, within the ADHD group, Node 9110, part of the cerebellum, appears densely connected to many regions of the cerebrum. Modern neuroscience, however, posits that the cerebellum and cerebrum serve relatively independent functions [18], which suggests that such extensive connections are unlikely to occur. In contrast, our method tends to yield a graph in which the node degrees are more evenly distributed, thereby offering a more balanced and potentially accurate representation.

## 7. Conclusion

We propose a neighborhood selection method for estimating the structure of a functional graphical model and show that it can consistently recover the conditional independence graph in the high-dimensional setting. Specifically, we pose the problem of graph selection as a series of function-on-function regressions, and we approximate the function-on-function regressions with a vector-on-vector regression approach that is achieved by functional dimension reduction. Through extensive simulations, we demonstrate that the proposed method outperforms existing approaches in a variety of settings. Finally, we apply our method on fMRI data sets that include patients with ASD and patients with ADHD, as well as corresponding control groups. We estimate the connectivity pattern between brain regions and find results that agree with previous neuroscience research.

A key step in our method is the choice of the basis for dimension reduction. Although we suggest using the FPCA basis for most settings, our methodology allows an arbitrary orthonormal basis. We also provide a theoretically motivated

procedure for choosing a particular basis. However, developing a more rigorous data-driven approach is still an open problem that we hope to study in the future. Another fruitful avenue for future work is the development of methods that allow for inference and hypothesis testing in functional graphs. For example, [26] has developed inferential tools for high-dimensional Markov networks, and future work may extend their results to the functional graph setting.

**Appendix A: Technical proofs**

We give proofs of the technical results that appear in the main text.

**A.1. Proof of Theorem 2.1**

For all  $k \in [p]$  and  $k \neq j$ , we define  $\mathcal{B}_{jk} : \mathbb{H} \mapsto \mathbb{H}$  as

$$\mathcal{B}_{jk}(h) := \mathcal{B}_j((0, \dots, 0, \underset{k\text{-th}}{h}, 0, \dots, 0)) \quad \text{for all } h \in \mathbb{H}.$$

Since  $\mathcal{B}_j$  is Hilbert-Schmidt, we claim that  $\mathcal{B}_{jk} \in \mathcal{B}_{\text{HS}}(\mathbb{H})$ . To prove this claim, note that for any CONS of  $\mathbb{H}$  denoted by  $\{e_n\}_{n=1}^\infty$ , we have

$$\{ \{(e_n, 0, \dots, 0)\}_{n=1}^\infty, \{(0, e_n, 0, \dots, 0)\}_{n=1}^\infty, \dots, \{(0, \dots, 0, e_n)\}_{n=1}^\infty \}$$

to be a CONS of  $\mathbb{H}^{p-1}$ . Given the assumption that  $\mathcal{B}_j$  is Hilbert-Schmidt, we have

$$\begin{aligned} \sum_{n=1}^\infty \|\mathcal{B}_{jk}(e_n)\|^2 &= \sum_{n=1}^\infty \left\| \mathcal{B}_j \left( (0, \dots, 0, \underset{k\text{-th}}{e_n}, \dots, 0) \right) \right\|^2 \\ &\leq \sum_{n=1}^\infty \|\mathcal{B}_j((e_n, 0, \dots, 0))\|^2 + \dots + \|\mathcal{B}_j((0, \dots, 0, e_n))\|^2 \\ &< \infty, \end{aligned}$$

which implies that  $\mathcal{B}_{jk} \in \mathcal{B}_{\text{HS}}(\mathbb{H})$ . By the linearity of  $\mathcal{B}_j$ , then for all

$$h = (h_1, \dots, h_{j-1}, h_{j+1}, \dots, h_p) \in \mathbb{H}^{p-1},$$

we have

$$\mathcal{B}_j(h) = \mathcal{B}_j((h_1, 0, \dots, 0)) + \dots + \mathcal{B}_j((0, \dots, 0, h_p)) = \sum_{k=1, k \neq j}^p \mathcal{B}_{jk}(h_k).$$

Thus, we have

$$\mathbb{E}[g_j \mid \mathbf{g}_{-j}] = \mathcal{B}_j(\mathbf{g}_{-j}) = \sum_{k \neq j} \mathcal{B}_{jk}(g_k). \tag{41}$$

The rest of the proof is composed of two steps. We first construct functions  $\{\beta_{jk}(t, t')\}_{k \neq j}$  such that (5) holds and then show (6) and (7). For any choice of

CONS  $\{\phi_m\}_{m=1}^\infty$  for  $\mathbb{H}$ , by Theorem 4.4.5 of [22] and the fact that  $\mathcal{B}_{jk} \in \mathcal{B}_{\text{HS}}(\mathbb{H})$ , we have

$$\mathcal{B}_{jk} = \sum_{m=1}^\infty \sum_{m'=1}^\infty b_{jk,mm'}^* \phi_m \otimes \phi_{m'},$$

where  $b_{jk,mm'}^* := \langle \mathcal{B}_{jk}(\phi_{m'}), \phi_m \rangle$  and  $\|\mathcal{B}_{jk}\|_{\text{HS}}^2 = \sum_{m=1}^\infty \sum_{m'=1}^\infty (b_{jk,mm'}^*)^2 < \infty$ . Let

$$\beta_{jk}(t, t') = \sum_{m, m'=1}^\infty b_{jk,mm'}^* \phi_m(t) \phi_{m'}(t')$$

for all  $(t, t') \in \mathcal{T} \times \mathcal{T}$ . Then, for any  $h \in \mathbb{H}$ , we have

$$\begin{aligned} \mathcal{B}_{jk}(h)(t) &= \sum_{m=1}^\infty \sum_{m'=1}^\infty b_{jk,mm'}^* \langle h, \phi_{m'} \rangle \phi_m(t) \\ &= \sum_{m=1}^\infty \sum_{m'=1}^\infty b_{jk,mm'}^* \left( \int_{\mathcal{T}} h(t') \phi_{m'}(t') dt' \right) \phi_m(t) \\ &= \int_{\mathcal{T}} \sum_{m=1}^\infty \sum_{m'=1}^\infty b_{jk,mm'}^* \phi_m(t) \phi_{m'}(t') h(t') dt' \\ &= \int_{\mathcal{T}} \beta_{jk}(t, t') h(t') dt' \end{aligned} \tag{42}$$

for all  $t \in \mathcal{T}$ , where the third equality is by Fubini's Theorem. In this way,  $\mathcal{B}_{jk}$  is the integral operator with the kernel  $\beta_{jk}(t, t')$ . By Theorem 4.6.7 of [22], we have  $\|\beta_{jk}\|_{\text{HS}} = \|\mathcal{B}_{jk}\|_{\text{HS}} < \infty$ . The relation (5) follows by combining (42) and (41).

We then show (6) and (7). Let  $\{\tilde{\phi}_m\}_{m=1}^\infty$  be another CONS of  $\mathbb{H}$ . Let  $\tilde{b}_{jk,mm'}^* := \langle \mathcal{B}_{jk}(\tilde{\phi}_{m'}), \tilde{\phi}_m \rangle$ , and we can similarly define  $\tilde{\beta}_{jk}(t, t')$  by

$$\tilde{\beta}_{jk}(t, t') = \sum_{m, m'=1}^\infty \tilde{b}_{jk,mm'}^* \tilde{\phi}_m(t) \tilde{\phi}_{m'}(t').$$

Similar to (42), we can show  $\mathcal{B}_{jk}(h)(t) = \int_{\mathcal{T}} \tilde{\beta}_{jk}(t, t') h(t') dt'$  for all  $h \in \mathbb{H}$  and  $t \in \mathcal{T}$ . Thus, we have

$$\tilde{b}_{jk,mm'}^* = \langle \mathcal{B}_{jk}(\tilde{\phi}_{m'}), \tilde{\phi}_m \rangle = \int_{\mathcal{T}} \tilde{\beta}_{jk}(t, t') \tilde{\phi}_m(t) \tilde{\phi}_{m'}(t') dt' dt.$$

In this way, to finish the proof, we only need to show that  $\beta_{jk}(t, t') = \tilde{\beta}_{jk}(t, t')$  a.e., or equivalently  $\|\beta_{jk} - \tilde{\beta}_{jk}\|_{\text{HS}} = 0$ . This is obvious since, for any CONS  $\{\bar{\phi}_m\}_{m \geq 1}$  for  $\mathbb{H}$ , we have

$$\mathcal{B}_{jk}(\bar{\phi}_m)(t) = \int_{\mathcal{T}} \beta_{jk}(t, t') \bar{\phi}_m(t') dt' = \int_{\mathcal{T}} \tilde{\beta}_{jk}(t, t') \bar{\phi}_m(t') dt'$$

for all  $t \in \mathcal{T}$ , which implies that

$$\|\beta_{jk} - \tilde{\beta}_{jk}\|_{\text{HS}}^2 = \sum_{m=1}^\infty \|\mathcal{B}_{jk}(\bar{\phi}_m) - \mathcal{B}_{jk}(\bar{\phi}_m)\|^2 = \sum_{m=1}^\infty 0 = 0.$$

**A.2. Derivation of (10) and (15)**

Recall that  $\beta_{jk}(t, t')$  is defined in (5) and  $\hat{\phi}_{jm}$  is an estimate of the true basis function  $\phi_{jm}$ . Let  $b_{jk,mm'}^* = \int_{\mathcal{T} \times \mathcal{T}} \beta_{jk}(t', t) \phi_m(t) \phi_{m'}(t') dt' dt$ . We focus on a given node  $j \in [p]$ , and we drop the index  $j$  from the notation to simplify the discussion.

We first prove (10). By (5) and (6), we have

$$\begin{aligned}
 a_{im}^Y &= \sum_{k=1}^{p-1} \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') g_i^{X_k}(t') \phi_m(t) dt' dt + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \int_{\mathcal{T} \times \mathcal{T}} \left( \sum_{m'', m'=1}^{\infty} b_{k, m'' m'}^* \phi_{m''}(t) \phi_{m'}(t') \right) g_i^{X_k}(t') \phi_m(t) dt' dt \\
 &\quad + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \sum_{m', m''=1}^{\infty} b_{k, m'' m'}^* \int_{\mathcal{T} \times \mathcal{T}} \phi_{m''}(t) \phi_{m'}(t') g_i^{X_k}(t') \phi_m(t) dt' dt \\
 &\quad + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \sum_{m', m''=1}^{\infty} b_{k, m'' m'}^* \left( \int_{\mathcal{T}} \phi_{m''}(t) \phi_m(t) dt \right) \left( \int_{\mathcal{T}} \phi_{m'}(t') g_i^{X_k}(t') dt' \right) \\
 &\quad + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \sum_{m'=1}^{\infty} b_{k, mm'}^* \left( \int_{\mathcal{T}} g_i^{X_k}(t') \phi_{m'}(t') dt' \right) + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \sum_{m'=1}^{\infty} b_{k, mm'}^* a_{im'}^{X_k} + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt \\
 &= \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} b_{k, mm'}^* a_{im'}^{X_k} + \sum_{k=1}^{p-1} \sum_{m'=1}^M b_{k, mm'}^* a_{im'}^{X_k} + \int_{\mathcal{T}} e_i(t) \phi_m(t) dt. \tag{43}
 \end{aligned}$$

Then (10) follows directly from (43) by setting  $\mathbf{B}_{k, \mathbf{M}}^* = (b_{k, mm'}^*)_{1 \leq m, m' \leq M}$ ,

$$r_{im} = \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} b_{k, mm'}^* a_{im'}^{X_k},$$

$\mathbf{r}_{i, \mathbf{M}} = (r_{i1}, \dots, r_{iM})^\top$ ,  $w_{im} = \int_{\mathcal{T}} e_i(t) \phi_m(t) dt$  and  $\mathbf{w}_{i, \mathbf{M}} = (w_{i1}, \dots, w_{iM})^\top$ .

To show (15), we only need to redefine relevant concepts. We define

$$\tilde{b}_{k, mm'} = \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t', t) \hat{\phi}_m(t') \hat{\phi}_{m'}(t) dt' dt. \tag{44}$$

Thus  $\tilde{b}_{k,mm'} = 0$  for all  $m, m' \geq 1$  when  $k \notin \mathcal{N}_j$ . Similarly, let

$$\tilde{w}_{im} = \int_{\mathcal{T}} e_i(t) \hat{\phi}_m(t) dt \quad \text{and} \quad \tilde{r}_{im} = \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} \tilde{b}_{k,mm'} \hat{a}_{im'}^{X_k}.$$

Furthermore, let  $\tilde{\mathbf{w}}_{i,M} = (\tilde{w}_{i1}, \dots, \tilde{w}_{iM})^\top$ ,  $\tilde{\mathbf{r}}_{i,M} = (\tilde{r}_{i1}, \dots, \tilde{r}_{iM})^\top \in \mathbb{R}^M$ ,  $\tilde{\mathbf{B}}_l^M$  is a  $M \times M$  matrix with the  $m$ -th row equal to  $(\tilde{b}_{l,m1}, \dots, \tilde{b}_{l,mM})$ , and

$$\mathbf{v}_{iM} = \sum_{k=1}^{p-1} (\tilde{\mathbf{B}}_{k,M} - \mathbf{B}_{k,M}^*) \hat{\mathbf{a}}_{i,M}^{X_k} + (\tilde{\mathbf{r}}_{i,M} - \mathbf{r}_{i,M}) + (\tilde{\mathbf{w}}_{i,M} - \mathbf{w}_{i,M}). \quad (45)$$

By (6), we have

$$\beta_k(t', t) = \sum_{m, m'=1}^{\infty} \tilde{b}_{k,mm'}^* \hat{\phi}_m(t) \hat{\phi}_{m'}(t') \quad \text{almost everywhere.}$$

Then by a similar argument to (43), we have

$$\hat{a}_{im}^Y = \sum_{k=1}^{p-1} \sum_{m' > M}^{\infty} \tilde{b}_{k,mm'} \hat{a}_{im'}^{X_k} + \sum_{k=1}^{p-1} \sum_{m'=1}^M \tilde{b}_{k,mm'} \hat{a}_{im'}^{X_k} + \int_{\mathcal{T}} e_i(t) \hat{\phi}_m(t) dt, \quad (46)$$

which implies (15) combined with (45).

### A.3. Simplification of ADMM optimization problems

We explain how to obtain the problem in (21). Let  $g(\sum_{k=1}^{p-1} \mathbf{Q}_k) = \frac{1}{2n} \|\mathbf{A}^Y - \sum_{k=1}^{p-1} \mathbf{Q}_k\|_F^2$  and  $\mathbf{W}_k^{h+1} = \mathbf{A}^{X_k} \mathbf{P}_k^{h+1} + \mathbf{U}_k^h$ . The update to matrices  $\{\mathbf{Q}_k\}_{k \in [p-1]}$  can be rewritten as

$$\min_{\{\mathbf{Q}_k\}_{k \in [p-1]}} g((p-1)\bar{\mathbf{Q}}) + \frac{\rho}{2} \sum_{k=1}^{p-1} \|\mathbf{Q}_k - \mathbf{W}_k^{h+1}\|_F^2 \quad \text{subject to} \quad \bar{\mathbf{Q}} = \frac{1}{p-1} \sum_{k=1}^{p-1} \mathbf{Q}_k.$$

Let  $\phi : \mathbb{R}^{n \times M} \rightarrow \mathbb{R}$  be a function that satisfies  $\phi(\mathbf{0}) = 0$  and  $\nabla \phi(\mathbf{0}) \neq 0$ . The Lagrangian function is then

$$\begin{aligned} & (\mathbf{Q}_1, \dots, \mathbf{Q}_{p-1}; \mu) = \\ & g((p-1)\bar{\mathbf{Q}}) + \frac{\rho}{2} \sum_{k=1}^{p-1} \|\mathbf{Q}_k - \mathbf{W}_k^{h+1}\|_F^2 + \mu \phi \left( \bar{\mathbf{Q}} - \frac{1}{p-1} \sum_{k=1}^{p-1} \mathbf{Q}_k \right). \end{aligned}$$

The matrix  $\mathbf{Q}_k$  that minimizes the Lagrangian satisfies

$$\frac{\partial L}{\partial \mathbf{Q}_k} = \rho(\mathbf{Q}_k - \mathbf{W}_k^{h+1}) - \frac{\mu}{p-1} \nabla \phi \left( \bar{\mathbf{Q}} - \frac{1}{p-1} \sum_{k=1}^{p-1} \mathbf{Q}_k \right) = 0,$$

$$\frac{\partial L}{\partial \mu} = \bar{Q} - \frac{1}{p-1} \sum_{k=1}^{p-1} Q_k = 0.$$

This is equivalent to  $\rho(Q_k - W_k^{h+1}) = \frac{\mu}{p-1} \nabla \phi(\mathbf{0})$ . As a result, we see that each entry in  $Q_k - W_k^{h+1}$  does not vary with  $k$ . Let  $\bar{R}^{h+1} = \frac{1}{p-1} \sum_{k=1}^{p-1} W_k^{h+1}$ . Then  $Q_k^{h+1}$  can be replaced by  $\bar{Q}^{h+1} + W_k^{h+1} - \bar{R}^{h+1}$  and

$$Q_k^{h+1} = \bar{Q}^{h+1} + A^{X_k} P_k^{h+1} + W_k^h - \overline{A^X P}^{h+1} - \bar{U}^h,$$

where  $\overline{A^X P}^h = \frac{1}{p-1} \sum_{k=1}^{p-1} A^{X_k} P_k^h$ . Therefore, we have obtained the problem in (21).

**A.4. Derivation of (40)**

We drop the subscript  $M$ . By (10) and the definition of  $\mathcal{N}_j$ , we have

$$a_i^Y = \sum_{k \in \mathcal{N}_j} B_k^* a_i^{X_k} + w_i + r_i,$$

where  $r_i = (r_{i1}, r_{i2}, \dots, r_{iM})^\top$  with

$$r_{im} = \sum_{k \in \mathcal{N}_j} \sum_{m' > M} b_{k,mm'}^* a_{im'}^{X_k}.$$

Let  $\Sigma^Y = \mathbb{E} [a_i^Y (a_i^Y)^\top]$ ,  $\Sigma^{Y, X_k} = \mathbb{E} [a_i^Y (a_i^{X_k})^\top]$ ,  $k \in [p-1]$ . Note that

$$((a_i^Y)^\top, (a_i^{X_1})^\top, (a_i^{X_2})^\top, \dots, (a_i^{X_{p-1}})^\top)^\top$$

is a multivariate Gaussian vector. Then

$$[B_k^*]_{k \in \mathcal{N}_j} = [\Sigma^{Y, X_k} - \Sigma^{r, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \tag{47}$$

and

$$\begin{aligned} & \left\| [B_k^*]_{k \in \mathcal{N}_j} \right\|_{\mathbb{F}} \\ &= \left\| [\Sigma^{Y, X_k} - \Sigma^{r, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \right\|_{\mathbb{F}} \\ &\geq \left\| [\Sigma^{Y, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \right\|_{\mathbb{F}} - \left\| [\Sigma^{r, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \right\|_{\mathbb{F}}. \end{aligned}$$

Since the correlation between  $a_i^Y$  and  $a_i^{X_k}$  is larger than correlation between  $r_i$  and  $a_i^{X_k}$ ,  $k \in \mathcal{N}_j$ , when  $M$  is large enough, we will have

$$\left\| [\Sigma^{Y, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \right\|_{\mathbb{F}} \gg \left\| [\Sigma^{r, X_k}]_{k \in \mathcal{N}_j} \left( \Sigma_{\mathcal{N}_j, \mathcal{N}_j}^X \right)^{-1} \right\|_{\mathbb{F}}.$$

Combining the last two displays, we have

$$\left\| [\mathbf{B}_{\mathbf{k}}^*]_{k \in \mathcal{N}_j} \right\|_{\mathbb{F}} \gtrsim \left\| [\boldsymbol{\Sigma}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \right\|_{\mathbb{F}}. \quad (48)$$

Furthermore, let  $\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k} = (\boldsymbol{\Sigma}^{\mathbf{Y}})^{-1/2} \boldsymbol{\Sigma}^{\mathbf{Y}, \mathbf{X}_k} (\boldsymbol{\Sigma}^{\mathbf{X}_k})^{-1/2}$ ,  $k \in [p-1]$ , and

$$\mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} = \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1/2} \right]_{k \in \mathcal{N}_j} \right) \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1/2} \right]_{k \in \mathcal{N}_j} \right).$$

Then

$$\begin{aligned} & [\boldsymbol{\Sigma}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \\ &= (\boldsymbol{\Sigma}^{\mathbf{Y}})^{\frac{1}{2}} [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{1/2} \right]_{k \in \mathcal{N}_j} \right) \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \\ &= (\boldsymbol{\Sigma}^{\mathbf{Y}})^{\frac{1}{2}} [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1/2} \right]_{k \in \mathcal{N}_j} \right) \end{aligned}$$

and

$$\begin{aligned} & [\boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} = \\ & (\boldsymbol{\Sigma}^{\mathbf{r}})^{\frac{1}{2}} [\mathbf{R}^{\mathbf{r}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1/2} \right]_{k \in \mathcal{N}_j} \right). \end{aligned}$$

By Lemma B.15, we have

$$\begin{aligned} & \left\| [\boldsymbol{\Sigma}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \right\|_{\mathbb{F}} \\ & \geq \left\| (\boldsymbol{\Sigma}^{\mathbf{Y}})^{\frac{1}{2}} \right\|_{\mathbb{F}} \left\{ \rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \right. \right. \\ & \quad \left. \left. \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1} \right]_{k \in \mathcal{N}_j} \right) \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right) \right\}^{1/2} \\ & = \sqrt{\text{tr}(\boldsymbol{\Sigma}^{\mathbf{Y}})} \left\{ \rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \right. \right. \\ & \quad \left. \left. \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1} \right]_{k \in \mathcal{N}_j} \right) \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right) \right\}^{1/2} \\ & \geq \sqrt{\text{tr}(\boldsymbol{\Sigma}^{\mathbf{Y}})} \sqrt{\rho_{\min} \left( \text{diag} \left( \left[ (\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})^{-1} \right]_{k \in \mathcal{N}_j} \right) \right)}. \\ & \quad \sqrt{\rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-2} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right)} \\ & \geq \sqrt{\text{tr}(\boldsymbol{\Sigma}^{\mathbf{Y}})} \sqrt{\frac{1}{\max_{k \in \mathcal{N}_j} \rho_{\max}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})}}}. \end{aligned}$$



$$\sqrt{\rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-2} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right)},$$

and, therefore,

$$\begin{aligned} & \sqrt{\kappa(M)} \left\| [\boldsymbol{\Sigma}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-1} \right\|_{\mathbb{F}} \\ & \geq \sqrt{\text{tr}(\boldsymbol{\Sigma}^{\mathbf{Y}})} \sqrt{\frac{\rho_{\min}(\boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}})}{\max_{k \in \mathcal{N}_j} \rho_{\max}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})}}} \\ & \sqrt{\rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-2} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right)} \\ & \geq \sqrt{\text{tr}(\boldsymbol{\Sigma}^{\mathbf{Y}})} \sqrt{\frac{\rho_{\min}(\boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}})}{\rho_{\max}(\boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}})}}} \\ & \sqrt{\rho_{\min} \left( [\mathbf{R}^{\mathbf{Y}, \mathbf{X}_k}]_{k \in \mathcal{N}_j} \left( \mathbf{R}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \right)^{-2} [\mathbf{R}^{\mathbf{X}_k, \mathbf{Y}}]_{k \in \mathcal{N}_j} \right)}. \end{aligned} \quad (49)$$

Combining (48), (49), and the definition of  $\Lambda(M, \phi)$ , we arrive at (40).

#### A.5. Proposition A.1 and its proof

**Proposition A.1.** *If the tuning parameter  $\lambda_n$  satisfies*

$$\lambda_n > \max_{k \in [p-1]} \frac{1}{n} \|(\mathbf{A}^{\mathbf{X}_k})^\top \mathbf{A}^{\mathbf{Y}}\|_F, \quad (50)$$

where  $\mathbf{A}^{\mathbf{X}_k}$  and  $\mathbf{A}^{\mathbf{Y}}$  are defined in Section 3, then the estimated support set  $\hat{\mathcal{N}}_j$  is empty.

*Proof.* This threshold of  $\lambda_n$  can be derived using the KKT condition. We use the notation introduced in Section 3. The subgradient of the objective in (13) with respect to  $\mathbf{B}_k$  can be written as

$$-\frac{1}{n} (\mathbf{A}^{\mathbf{X}_k})^\top \left( \mathbf{A}^{\mathbf{Y}} - \sum_{l=1}^{p-1} \mathbf{A}^{\mathbf{X}_l} \mathbf{B}_l \right) + \lambda_n \boldsymbol{\Upsilon}_k, \quad (51)$$

where  $\boldsymbol{\Upsilon}_k = \|\mathbf{B}_k\|_F^{-1} \mathbf{B}_k$  if  $\mathbf{B}_k \neq \mathbf{0}$ , and  $\boldsymbol{\Upsilon}_k \in \mathbb{R}^{M \times M}$ ,  $\|\boldsymbol{\Upsilon}_k\|_F \leq 1$  otherwise. We assume that  $\hat{\mathcal{N}}_j$  is non-empty. That is, there exists some  $k$  such that  $\mathbf{B}_k \neq \mathbf{0}$ . By (51), we have

$$\lambda_n \frac{\mathbf{B}_k}{\|\mathbf{B}_k\|_F} = \frac{1}{n} (\mathbf{A}^{\mathbf{X}_k})^\top \mathbf{A}^{\mathbf{Y}} - \frac{1}{n} (\mathbf{A}^{\mathbf{X}_k})^\top \sum_{l \in \hat{\mathcal{N}}_j} \mathbf{A}^{\mathbf{X}_l} \mathbf{B}_l \quad \text{for all } k \in \hat{\mathcal{N}}_j. \quad (52)$$

Let  $\mathbf{A}^{\mathcal{X}_{\hat{\mathcal{N}}_j}} = [\mathbf{A}^{\mathbf{X}_k}]_{k \in \hat{\mathcal{N}}_j} \in \mathbb{R}^{n \times |\hat{\mathcal{N}}_j| M}$  and  $\mathbf{B}^{\mathcal{X}_{\hat{\mathcal{N}}_j}} = [\mathbf{B}_l]_{l \in \hat{\mathcal{N}}_j} \in \mathbb{R}^{M \times |\hat{\mathcal{N}}_j| M}$  be the submatrix of  $\mathbf{A}^{\mathbf{X}}$  and  $[\mathbf{B}_1^\top, \dots, \mathbf{B}_p^\top]^\top$  that correspond to  $k \in \hat{\mathcal{N}}_j$ . By (52), we have

$$\left( \frac{1}{n} (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})^\top (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}}) + \lambda_n \text{diag} \left( \left\{ \frac{1}{\|\mathbf{B}_l\|_F} \mathbf{I}_M \right\}_{l \in \hat{\mathcal{N}}_j} \right) \right) \mathbf{B}^{\mathcal{X}_{\hat{\mathcal{N}}_j}} = \frac{1}{n} (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})^\top \mathbf{A}^{\mathbf{Y}}. \quad (53)$$

Since  $\frac{1}{n} (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})^\top (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})$  is positive semi-definite and we have assumed that (50) holds, the left hand side of (53) then satisfies that

$$\begin{aligned} & \left\| \left( \frac{1}{n} (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})^\top (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}}) + \lambda_n \text{diag} \left( \left\{ \frac{1}{\|\mathbf{B}_l\|_F} \mathbf{I}_M \right\}_{l \in \hat{\mathcal{N}}_j} \right) \right) \mathbf{B}^{\mathcal{X}_{\hat{\mathcal{N}}_j}} \right\|_F \\ & \geq \left\| \lambda_n \text{diag} \left( \left\{ \frac{1}{\|\mathbf{B}_l\|_F} \mathbf{I}_M \right\}_{l \in \hat{\mathcal{N}}_j} \right) \mathbf{B}^{\mathcal{X}_{\hat{\mathcal{N}}_j}} \right\|_F \\ & \geq \lambda_n |\hat{\mathcal{N}}_j| \\ & > |\hat{\mathcal{N}}_j| \cdot \max_{k \in [p-1]} \frac{1}{n} \|(\mathbf{A}^{\mathbf{X}_k})^\top \mathbf{A}^{\mathbf{Y}}\|_F, \end{aligned}$$

where the first inequality follows from Lemma B.16. On the other hand, the right hand side of (53) satisfies that

$$\left\| \frac{1}{n} (\mathbf{A}^{\mathbf{X}_{\hat{\mathcal{N}}_j}})^\top \mathbf{A}^{\mathbf{Y}} \right\|_F \leq |\hat{\mathcal{N}}_j| \cdot \max_{k \in [p-1]} \frac{1}{n} \|(\mathbf{A}^{\mathbf{X}_k})^\top \mathbf{A}^{\mathbf{Y}}\|_F.$$

Combine the above two equations with (53), we have a contradiction. Thus, we conclude that  $\hat{\mathcal{N}}_j$  must be empty.  $\square$

#### A.6. Proof of Theorem 4.1

In this section, we prove Theorem 4.1. We first introduce some useful notation.

Let  $\tilde{\lambda}(n, p, M, \delta)$  be defined as

$$\begin{aligned} \tilde{\lambda}(n, p, M, \delta) = & 2\mathcal{C}_{n, \delta} \left( \frac{M \sqrt{\Xi_1(M)}}{\sqrt{n}} + 2\sqrt{\Xi_1(M)} \sqrt{\frac{\log(4(p-1)/\delta)}{n}} \right) \\ & + 2\omega(M) \left\{ 7\sqrt{3} \frac{M \sqrt{\log(6(p-1)/\delta) + 2 \log M}}{\sqrt{n}} \right. \\ & \left. + \frac{8Mc(\log(2n))(\log(6(p-1)/\delta) + 2 \log M)}{3n} \right\}, \quad (54) \end{aligned}$$

where  $c$  is some universal constant that does not depend on  $n$ ,  $p$  or  $M$ .

To simplify the notation, we omit the basis dimension,  $M$ , and let  $\mathbf{a}_i^Y = \mathbf{a}_{i,M}^Y$ ,  $\mathbf{a}_i^{X_k} = \mathbf{a}_{i,M}^{X_k}$ , and  $\mathbf{B}_k^* = \mathbf{B}_{k,M}^*$  for all  $k \in [p-1]$ . Then by (10), for all  $i \in [n]$ , we have

$$\mathbf{a}_i^Y = \sum_{k=1}^{p-1} \mathbf{B}_k^* \mathbf{a}_i^{X_k} + \mathbf{u}_i, \quad (55)$$

where  $\mathbf{u}_i = \mathbf{w}_i + \mathbf{r}_i$ , and  $\mathbf{w}_i, \mathbf{r}_i$  are defined in Appendix A.2. With this notation, we give a proof of Theorem 4.1.

*Proof.* The equation (55) can be rewritten as

$$\mathbf{a}_i^Y = \sum_{k=1}^{p-1} \left( (\mathbf{a}_i^{X_k})^\top \otimes \mathbf{I}_M \right) \text{vec}(\mathbf{B}_k^*) + \mathbf{u}_i. \quad (56)$$

Let  $\mathbf{Z}_i^{X_k} = \mathbf{a}_i^{X_k} \otimes \mathbf{I}_M \in \mathbb{R}^{M^2 \times M}$ ,  $i \in [n]$ , and let  $\beta_k^* = \text{vec}(\mathbf{B}_k^*) \in \mathbb{R}^{M^2}$ ,  $k \in [p-1]$ . Furthermore, let

$$\begin{aligned} \mathbf{a}^Y &= ((\mathbf{a}_1^Y)^\top, (\mathbf{a}_2^Y)^\top, \dots, (\mathbf{a}_n^Y)^\top)^\top \in \mathbb{R}^{nM}, \\ \beta^* &= ((\beta_1^*)^\top, (\beta_2^*)^\top, \dots, (\beta_{p-1}^*)^\top)^\top \in \mathbb{R}^{(p-1)M^2}, \end{aligned}$$

and

$$\mathbf{u} = (\mathbf{u}_1^\top, \mathbf{u}_2^\top, \dots, \mathbf{u}_n^\top)^\top \in \mathbb{R}^{nM}. \quad (57)$$

Let  $\mathbf{Z}_i = ((\mathbf{Z}_i^{X_1})^\top, (\mathbf{Z}_i^{X_2})^\top, \dots, (\mathbf{Z}_i^{X_{p-1}})^\top)^\top \in \mathbb{R}^{(p-1)M^2 \times M}$  for  $i \in [n]$ , and

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1^\top \\ \mathbf{Z}_2^\top \\ \vdots \\ \mathbf{Z}_n^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{Z}_1^{X_1})^\top & (\mathbf{Z}_1^{X_2})^\top & \dots & (\mathbf{Z}_1^{X_{p-1}})^\top \\ (\mathbf{Z}_2^{X_1})^\top & (\mathbf{Z}_2^{X_2})^\top & \dots & (\mathbf{Z}_2^{X_{p-1}})^\top \\ \vdots & \vdots & \dots & \vdots \\ (\mathbf{Z}_n^{X_1})^\top & (\mathbf{Z}_n^{X_2})^\top & \dots & (\mathbf{Z}_n^{X_{p-1}})^\top \end{bmatrix} \in \mathbb{R}^{nM \times (p-1)M^2}.$$

Then we can further formulate (56) as

$$\mathbf{a}^Y = \mathbf{Z} \beta^* + \mathbf{u}. \quad (58)$$

Recall that

$$\mathbf{A}^X = \begin{bmatrix} (\mathbf{a}_1^X)^\top \\ (\mathbf{a}_2^X)^\top \\ \vdots \\ (\mathbf{a}_n^X)^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{a}_1^{X_1})^\top & (\mathbf{a}_1^{X_2})^\top & \dots & (\mathbf{a}_1^{X_{p-1}})^\top \\ (\mathbf{a}_2^{X_1})^\top & (\mathbf{a}_2^{X_2})^\top & \dots & (\mathbf{a}_2^{X_{p-1}})^\top \\ \vdots & \vdots & \dots & \vdots \\ (\mathbf{a}_n^{X_1})^\top & (\mathbf{a}_n^{X_2})^\top & \dots & (\mathbf{a}_n^{X_{p-1}})^\top \end{bmatrix} \in \mathbb{R}^{n \times (p-1)M},$$

then we have  $\mathbf{Z} = \mathbf{A}^X \otimes \mathbf{I}_M$ . We divide columns of  $\mathbf{A}^X$  into  $p-1$  groups with equal group size  $M$ , that is,  $\mathbf{A}^X = (\mathbf{A}^{X_1}, \mathbf{A}^{X_2}, \dots, \mathbf{A}^{X_{p-1}})$ , where  $\mathbf{A}^{X_k} \in \mathbb{R}^{n \times M}$  for all  $k \in [p-1]$ . Similarly, we divide the columns of  $\mathbf{Z}$  into  $p-1$  groups with equal group size  $M^2$ , that is,  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_{p-1})$ , where  $\mathbf{Z}_k \in$

$\mathbb{R}^{nM \times M^2}$  for all  $k \in [p-1]$ . Then, we have  $(\mathbf{Z}^\top \mathbf{u})_k = \mathbf{Z}_k^\top \mathbf{u}$ . Besides, by definition of  $\mathbf{Z}$ , it is easy to see that  $\mathbf{Z}_k = \mathbf{A}^{\mathbf{X}_k} \otimes \mathbf{I}_M$ .

Besides, we can rewrite (13) as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \{ \mathcal{L}(\boldsymbol{\beta}) + \lambda_n \mathcal{R}(\boldsymbol{\beta}) \}, \tag{59}$$

where

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{2n} \|\mathbf{a}^{\mathbf{Y}} - \mathbf{Z}\boldsymbol{\beta}\|_2^2, \tag{60}$$

$$\mathcal{R}(\boldsymbol{\beta}) = \sum_{k=1}^{p-1} \|\boldsymbol{\beta}_k\|_2, \tag{61}$$

and  $\boldsymbol{\beta} = ((\boldsymbol{\beta}_1)^\top, \dots, (\boldsymbol{\beta}_{p-1})^\top)^\top \in \mathbb{R}^{(p-1)M^2}$ , with  $\boldsymbol{\beta}_k \in \mathbb{R}^{M^2}$  for  $k \in [p-1]$ .

Thus, the support set defined in (8) and its estimator defined in (14) can be expressed as

$$\mathcal{N}_j = \{k \in [p-1] : \|\boldsymbol{\beta}_k^*\|_2 > 0\},$$

and

$$\hat{\mathcal{N}}_j = \left\{ k \in [p-1] : \|\hat{\boldsymbol{\beta}}_k\|_2 > \epsilon_n \right\}.$$

We define the model space  $\mathcal{M}(\mathcal{N}_j)$  with which the penalty term  $\mathcal{R}(\cdot)$  is decomposable. Let

$$\begin{aligned} \mathcal{M} &= \mathcal{M}(\mathcal{N}_j) = \\ &\{ \boldsymbol{\beta} = ((\boldsymbol{\beta}_1)^\top, \dots, (\boldsymbol{\beta}_{p-1})^\top)^\top \in \mathbb{R}^{(p-1)M^2} : \boldsymbol{\beta}_k = 0 \text{ for all } k \notin \mathcal{N}_j \}, \end{aligned}$$

we then have its orthogonal complement as

$$\begin{aligned} \mathcal{M}^\perp &= \mathcal{M}(\mathcal{N}_j)^\perp = \\ &\{ \boldsymbol{\beta} = ((\boldsymbol{\beta}_1)^\top, \dots, (\boldsymbol{\beta}_{p-1})^\top)^\top \in \mathbb{R}^{(p-1)M^2} : \boldsymbol{\beta}_k = 0 \text{ for all } k \in \mathcal{N}_j \}. \end{aligned}$$

It is then easy to verify that  $\mathcal{R}(\cdot)$  defined in (61) is decomposable with respect to  $(\mathcal{M}(\mathcal{N}_j), \mathcal{M}(\mathcal{N}_j)^\perp)$  (Example 2, Section 2.2 of [47]), that is

$$\mathcal{R}(\boldsymbol{\theta} + \boldsymbol{\gamma}) = \mathcal{R}(\boldsymbol{\theta}) + \mathcal{R}(\boldsymbol{\gamma}) \quad \text{for all } \boldsymbol{\theta} \in \mathcal{M}(\mathcal{N}_j) \text{ and } \boldsymbol{\gamma} \in \mathcal{M}(\mathcal{N}_j)^\perp.$$

When  $\lambda_n = \tilde{\lambda}(n, p, M, \delta)$ , where  $\tilde{\lambda}(n, p, M, \delta)$  is defined in (54), then by Lemma B.3, we have

$$\lambda_n \geq \frac{2}{n} \max_{k \in [p-1]} \|(\mathbf{Z}^\top \mathbf{u})_k\|_2 \tag{62}$$

hold with probability at least  $1 - 2\delta$ . This way, by Lemma 1 of [47], we have that  $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*$  lies in a constrained space  $\mathbb{C}(\mathcal{N}_j)$  defined by

$$\mathbb{C}(\mathcal{N}_j) = \{ \boldsymbol{\theta} \in \mathbb{R}^{(p-1)M^2} : \|\boldsymbol{\theta}_{\mathcal{M}^\perp}\|_{1,2} \leq 3\|\boldsymbol{\theta}_{\mathcal{M}}\|_{1,2} \}. \tag{63}$$

with probability at least  $1 - 2\delta$ . Note that  $\mathbb{C}(\mathcal{N}_j)$  depends on support  $\mathcal{N}_j$  through  $\mathcal{M}$ .

The error term of a first-order Taylor series expansion is

$$\mathcal{L}(\beta^* + \Delta\hat{\beta}) - \mathcal{L}(\beta^*) - \langle \nabla \mathcal{L}(\beta^*), \Delta\hat{\beta} \rangle = \frac{1}{2n} \|\mathbf{Z}\Delta\hat{\beta}\|_2^2$$

where  $\Delta\hat{\beta} = \hat{\beta} - \beta^*$ . Then by Lemma B.7 and  $\Delta\hat{\beta} \in \mathbb{C}(\mathcal{N}_j)$  with probability at least  $1 - 2\delta$ , we have

$$\mathbb{P} \left\{ \frac{1}{2n} \|\mathbf{Z}\Delta\hat{\beta}\|_2^2 \geq \frac{\kappa}{4} \|\Delta\hat{\beta}\|_2^2 \right\} \geq 1 - 3\delta. \tag{64}$$

Thus, by Lemma B.4, we then have

$$\mathbb{P} \left\{ \sqrt{\sum_{k=1}^{p-1} \|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F^2} \leq \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}} \right\} \geq 1 - 3\delta, \tag{65}$$

where  $\kappa(M)$  is defined in (28) and  $\chi(n, p, M, \delta)$  is defined in (32). Given the inequality in the left hand side parenthesis of (65) holds, we then have

$$\|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F \leq \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}} \quad \text{for all } k \in [p - 1].$$

Next, we prove that  $\mathbb{P}\{\hat{\mathcal{N}}_j = \mathcal{N}_j\} \geq 1 - 3\delta$ . To show that, we only need to prove that the above inequality implies that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$ . Under the assumption that the above inequality holds, note that for any  $k \notin \mathcal{N}_j$ , we have  $\mathbf{B}_k^* = \mathbf{0}$  for all  $M$ , thus, we have

$$\|\hat{\mathbf{B}}_k\|_F = \|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F \leq \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}}.$$

On the other hand, for  $k \in \mathcal{N}_j$ , we have

$$\begin{aligned} \|\hat{\mathbf{B}}_k\|_F &= \|\mathbf{B}_k^* + \hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F \geq \|\mathbf{B}_k^*\|_F - \|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F \\ &\geq \tau(M) - \chi(n, p, M, \delta) - 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}}. \end{aligned}$$

Because  $M \geq M^*$ , and by the definition of  $M^*$  in (31) and the definition of  $\nu(M)$  in (30), when

$$\chi(n, p, M, \delta) \leq \frac{\nu(M)}{3}, \tag{66}$$

we have

$$\tau(M) - \chi(n, p, M, \delta) - 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}} > \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}}.$$

Recall that  $\epsilon_n = \chi(n, p, M, \delta) + 12\sqrt{s/\kappa(M)}\omega(M)$ , we then have

$$\max_{k \notin \mathcal{N}_j} \|\hat{\mathbf{B}}_k\|_{\text{F}} \leq \epsilon_n < \min_{k \in \mathcal{N}_j} \|\hat{\mathbf{B}}_k\|_{\text{F}}.$$

which implies that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$  by (14).

Finally, we only need to show an asymptotic lower bound for  $n$  such that  $\Gamma(n, p, M, \delta) \leq \kappa(M)/(32M^2s)$  and  $\chi(n, p, M, \delta) \leq \nu(M)/3$  are both satisfied.

First, to satisfy  $\Gamma(n, p, M, \delta) \leq \kappa(M)/(32M^2s)$ , where  $\Gamma(n, p, M, \delta)$  is defined in (88), by dropping the  $\log n$  term, we have

$$n \geq \tilde{O}\left(\frac{M^4 s^2 \log(p^2 M^2 / \delta)}{\kappa^2(M)}\right).$$

Note that  $\Xi_i(M), i = 1, 2, 3$  are all uniformly bounded for all  $M$ . Next, to satisfy

$$\chi(n, p, M, \delta) \leq \nu(M)/3,$$

where  $\chi(n, p, M, \delta)$  is defined in (32), we need

$$n \geq \tilde{O}\left(\frac{s \cdot \max\{M^2, \log(p/\delta), M^2\omega^2(M) \log(M^2 p/\delta)\}}{\kappa(M)\nu^2(M)}\right).$$

Combine the above results and note that decreasing  $3\delta$  to  $\delta$  doesn't affect the asymptotic order of  $n$ , we then have the final result.  $\square$

### A.7. Proof of Theorem 4.3

In this section, we prove Theorem 4.3. In addition to the notations introduced in Appendix A.6, we define more notations that will be used in this section. Let

$$\begin{aligned} \mathcal{Q}_{n,\delta}^1 &:= 1 + 8 \left( \frac{\log(2/\delta)}{n} + \sqrt{\frac{\log(2/\delta)}{n}} \right), \\ \mathcal{Q}_{n,p,\delta}^2 &:= 1 + 8 \left( \frac{\log(2p/\delta)}{n} + \sqrt{\frac{\log(2p/\delta)}{n}} \right), \\ \mathcal{Q}_{n,p,M,\delta}^3 &:= 1 + 8 \left( \frac{\log(2(p-1)M/\delta)}{n} + \sqrt{\frac{\log(2(p-1)M/\delta)}{n}} \right). \end{aligned} \tag{67}$$

Also, let

$$\begin{aligned} \Gamma_v &= \frac{12\sigma_{\max,0}|\mathcal{N}_j|d_{j0}^2 \log(2p/\delta)}{n} \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{\text{HS}}^2 \right) \mathcal{Q}_{n,\delta}^1 \\ &\quad + 3\Phi(M)\sigma_{\max,0}|\mathcal{N}_j| \mathcal{Q}_{n,\delta}^1 (d_{j0}^2 - d_{js}^2(M)) \end{aligned}$$

$$+ 3\sigma_{jr}\mathcal{Q}_{n,\delta}^1\sqrt{d_{js}(M)}. \quad (68)$$

We then state the exact form of  $\check{\lambda}(n, p, M, \delta)$  as below:

$$\begin{aligned} \check{\lambda}(n, p, M, \delta) &= \check{\lambda}(n, p, M, \delta) + 2\Gamma_v^{1/2}\sqrt{\max_{k \in [p-1]} \text{tr}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})}\sqrt{\mathcal{Q}_{n,p-1,\delta}^2} \\ &+ \sqrt{\sigma_{\max,0}}\sqrt{\mathcal{Q}_{n,p,\delta}^2}\sqrt{\frac{(1/c_1)\log(c_2/\delta)}{n}}d_{js}(M) \cdot \left(\sqrt{\Xi_4(M)}\sqrt{\mathcal{Q}_{n,p,\delta}^2} + \Gamma_v^{1/2}\right). \end{aligned} \quad (69)$$

We also define that

$$\begin{aligned} \check{\Gamma}(n, p, M, \delta) &= \Gamma(n, p, M, \delta) + d_{j,\max}^2\sigma_{\max,0}\mathcal{Q}_{n,p,\delta}^2\frac{(1/c_1)\log(c_2/\delta)}{n} \\ &+ 2\sqrt{\Xi_2(M)}d_{j,\max}\sqrt{\sigma_{\max,0}}\sqrt{\mathcal{Q}_{n,p,\delta}^2}\sqrt{\frac{(1/c_1)\log(c_2/\delta)}{n}} \cdot \sqrt{\mathcal{Q}_{n,p,M,\delta}^3}. \end{aligned} \quad (70)$$

As in the proof of Theorem 4.1, we omit the basis dimension,  $M$ , and let  $\hat{\mathbf{a}}_i^{\mathbf{Y}} = \hat{\mathbf{a}}_{i,M}^{\mathbf{Y}}$ ,  $\hat{\mathbf{a}}_i^{\mathbf{X}_k} = \hat{\mathbf{a}}_{i,M}^{\mathbf{X}_k}$ , and  $\mathbf{B}_k^* = \mathbf{B}_{k,M}^*$  for all  $k \in [p-1]$ . Then by (15), for all  $i \in [n]$ , we have

$$\hat{\mathbf{a}}_i^{\mathbf{Y}} = \sum_{k=1}^{p-1} \mathbf{B}_k^* \hat{\mathbf{a}}_i^{\mathbf{X}_k} + \tilde{\mathbf{u}}_i, \quad (71)$$

where  $\tilde{\mathbf{u}}_i = \mathbf{w}_i + \mathbf{r}_i + \mathbf{v}_i$ , and  $\mathbf{w}_i, \mathbf{r}_i, \mathbf{v}_i$  are defined in Appendix A.2. We now give the proof of Theorem 4.3.

*Proof.* Let

$$\hat{\mathbf{A}}^{\mathbf{X}} = \begin{bmatrix} (\hat{\mathbf{a}}_1^{\mathbf{X}_1})^\top & (\hat{\mathbf{a}}_1^{\mathbf{X}_2})^\top & \cdots & (\hat{\mathbf{a}}_1^{\mathbf{X}_{p-1}})^\top \\ (\hat{\mathbf{a}}_2^{\mathbf{X}_1})^\top & (\hat{\mathbf{a}}_2^{\mathbf{X}_2})^\top & \cdots & (\hat{\mathbf{a}}_2^{\mathbf{X}_{p-1}})^\top \\ \vdots & \vdots & \ddots & \vdots \\ (\hat{\mathbf{a}}_n^{\mathbf{X}_1})^\top & (\hat{\mathbf{a}}_n^{\mathbf{X}_2})^\top & \cdots & (\hat{\mathbf{a}}_n^{\mathbf{X}_{p-1}})^\top \end{bmatrix} \in \mathbb{R}^{n \times (p-1)M}.$$

We divide columns of  $\hat{\mathbf{A}}^{\mathbf{X}}$  into  $p-1$  groups with equal group size  $M$ , that is,

$$\hat{\mathbf{A}}^{\mathbf{X}} = [\hat{\mathbf{A}}^{\mathbf{X}_1}, \hat{\mathbf{A}}^{\mathbf{X}_2}, \dots, \hat{\mathbf{A}}^{\mathbf{X}_{p-1}}],$$

where  $\hat{\mathbf{A}}^{\mathbf{X}_k} \in \mathbb{R}^{n \times M}$  for all  $k \in [p-1]$ . Let  $\hat{\mathbf{Z}} = \hat{\mathbf{A}}^{\mathbf{X}} \otimes \mathbf{I}_M$ . Similarly, we divide the columns of  $\hat{\mathbf{Z}}$  into  $p-1$  groups with equal group size  $M^2$ , that is,  $\hat{\mathbf{Z}}^{\mathbf{X}} = [\hat{\mathbf{Z}}^{\mathbf{X}_1}, \hat{\mathbf{Z}}^{\mathbf{X}_2}, \dots, \hat{\mathbf{Z}}^{\mathbf{X}_{p-1}}]$ , where  $\hat{\mathbf{Z}}^{\mathbf{X}_k} \in \mathbb{R}^{nM \times M^2}$  for all  $k \in [p-1]$ . Then, we have  $(\hat{\mathbf{Z}}^\top \mathbf{u})_k = \hat{\mathbf{Z}}_k^\top \mathbf{u}$ . By definition of  $\hat{\mathbf{Z}}$ , it is easy to see that  $\hat{\mathbf{Z}}_k = \hat{\mathbf{A}}^{\mathbf{X}_k} \otimes \mathbf{I}_M$ . In addition, let  $\hat{\mathbf{a}}^{\mathbf{Y}} = ((\hat{\mathbf{a}}_1^{\mathbf{Y}})^\top, \dots, (\hat{\mathbf{a}}_n^{\mathbf{Y}})^\top)^\top \in \mathbb{R}^{nM}$ . Furthermore, let  $\tilde{\boldsymbol{\Sigma}}_n^{\mathbf{X}} = \frac{1}{n}(\hat{\mathbf{A}}^{\mathbf{X}})^\top \hat{\mathbf{A}}^{\mathbf{X}}$ . Thus, we can rewrite (13) as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \hat{\mathcal{L}}(\boldsymbol{\beta}) + \lambda_n \mathcal{R}(\boldsymbol{\beta}) \right\},$$

where

$$\hat{\mathcal{L}}(\boldsymbol{\beta}) = \frac{1}{2n} \|\hat{\mathbf{a}}^Y - \hat{\mathbf{Z}}\boldsymbol{\beta}\|_2^2,$$

$$\mathcal{R}(\boldsymbol{\beta}) = \sum_{k=1}^{p-1} \|\boldsymbol{\beta}_k\|_2,$$

and  $\boldsymbol{\beta} = ((\boldsymbol{\beta}_1)^\top, \dots, (\boldsymbol{\beta}_{p-1})^\top)^\top \in \mathbb{R}^{(p-1)M^2}$ , with  $\boldsymbol{\beta}_k \in \mathbb{R}^{M^2}$  for  $k \in [p-1]$ .

We can then follow the similar proof of Theorem 4.1 to prove Theorem 4.3. The only two modifications needed are new upper bounds for

$$\frac{2}{n} \max_{k \in [p-1]} \|(\hat{\mathbf{Z}}^\top \tilde{\mathbf{u}})_k\|_2 \text{ and } \|\tilde{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty.$$

In fact, when we have

$$\lambda_n = \check{\lambda}(n, p, M, \delta) \geq \frac{2}{n} \max_{k \in [p-1]} \|(\hat{\mathbf{Z}}^\top \tilde{\mathbf{u}})_k\|_2 \text{ and } \|\tilde{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty \leq \check{\Gamma}(n, p, M, \delta) \tag{72}$$

hold, where  $\check{\lambda}(n, p, M, \delta)$  is defined in (38), then by the similar argument in the proof of Theorem 4.1, we can show that

$$\sqrt{\sum_{k=1}^{p-1} \|\hat{\mathbf{B}}_k - \mathbf{B}_k^*\|_F^2} \leq \check{\chi}(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}},$$

which further implies that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$ . Thus, to show that  $\hat{\mathcal{N}}_j = \mathcal{N}_j$  holds with high probability, we only need to prove that (72) holds with high probability.

We define the following events.

$$\mathcal{A}_1 = \left\{ \frac{1}{n} \sum_{i=1}^n \|g_j^{(i)}\|^2 \leq \sigma_{\max,0} \mathcal{Q}_{n,p,\delta}^2 \text{ for all } j \in [p] \right\},$$

$$\mathcal{A}_2 = \left\{ \frac{1}{n} \sum_{i=1}^n \|e^{(i)}\|^2 \leq \sigma_{jr} \mathcal{Q}_{n,p,\delta}^1 \right\},$$

$$\mathcal{A}_3 = \left\{ \frac{1}{n} \sum_{i=1}^n \|\mathbf{r}_i + \mathbf{w}_i\|_2^2 \leq \Xi_4(M) \mathcal{Q}_{n,p,\delta}^1 \right\},$$

$$\mathcal{A}_4 = \left\{ \frac{1}{n} \sum_{i=1}^n \|\mathbf{a}_i^{\mathbf{X}^k}\|_2^2 \leq \text{tr}(\boldsymbol{\Sigma}_{kk}^{\mathbf{X}}) \mathcal{Q}_{n,p-1,\delta}^2 \text{ for all } k \in [p-1] \right\},$$

$$\mathcal{A}_5 = \left\{ \frac{2}{n} \max_{k \in [p-1]} \|(\mathbf{Z}^\top (\mathbf{w} + \mathbf{r}))_k\|_2 \leq \tilde{\lambda}(n, p, M, \delta) \right\},$$

$$\mathcal{A}_6 = \left\{ \|\tilde{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty \leq \Gamma(n, p, M, \delta) \right\},$$

$$\mathcal{A}_7 = \left\{ \frac{1}{n} \sum_{i=1}^n \|a_{im}^{\mathbf{X}^k}\|_2^2 \leq \Xi_2(M) \mathcal{Q}_{n,p,\delta}^3 \text{ for all } m \in [M] \text{ and } k \in [p-1] \right\},$$



where  $\check{\lambda}(n, p, M, \delta)$  is defined in (33). Then we claim that under Assumption 4.4, we have

$$\begin{aligned} \cap_{i=1}^7 \mathcal{A}_i \implies \\ \frac{2}{n} \max_{k \in [p-1]} \|(\hat{\mathbf{Z}}^\top \tilde{\mathbf{u}})_k\|_2 \leq \check{\lambda}(n, p, M, \delta) \text{ and } \|\tilde{\Sigma}_n^{\mathbf{X}} - \Sigma^{\mathbf{X}}\|_\infty \leq \check{\Gamma}(n, p, M, \delta). \end{aligned} \quad (73)$$

We now prove the above claim. We first prove that  $\cap_{i=1}^7 \mathcal{A}_i$  implies that

$$\frac{2}{n} \max_{k \in [p-1]} \|(\hat{\mathbf{Z}}^\top \tilde{\mathbf{u}})_k\|_2 \leq \check{\lambda}(n, p, M, \delta).$$

Note that for all  $k \in [p-1]$ , we have

$$\begin{aligned} \|(\hat{\mathbf{Z}}^\top \tilde{\mathbf{u}})_k\|_2 &= \|\hat{\mathbf{Z}}_k^\top \tilde{\mathbf{u}}\|_2 \leq \|(\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}}\|_2 + \|\mathbf{Z}_k^\top \tilde{\mathbf{u}}\|_2 \\ &= \|(\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}}\|_2 + \|\mathbf{Z}_k^\top (\mathbf{w} + \mathbf{r} + \mathbf{v})\|_2 \\ &\leq \|(\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}}\|_2 + \|\mathbf{Z}_k^\top (\mathbf{w} + \mathbf{r})\|_2 + \|\mathbf{Z}_k^\top \mathbf{v}\|_2, \end{aligned}$$

where  $\mathbf{v} = (v_{im})_{\{m \in [M]\}}$  with  $v_{im}$  defined in (45). The above inequality implies that

$$\begin{aligned} \frac{2}{n} \max_{k \in [p-1]} \|(\mathbf{Z}^\top \mathbf{u})_k\|_2 &\leq \frac{2}{n} \max_{k \in [p-1]} \|(\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}}\|_2 \\ &\quad + \frac{2}{n} \max_{k \in [p-1]} \|\mathbf{Z}_k^\top (\mathbf{w} + \mathbf{r})\|_2 + \frac{2}{n} \max_{k \in [p-1]} \|\mathbf{Z}_k^\top \mathbf{v}\|_2 \end{aligned} \quad (74)$$

By  $\mathcal{A}_5$ , we have the second term to be bounded by  $\check{\lambda}(n, p, M, \delta)$ . To bound the third term, note that

$$\begin{aligned} \left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{v} \right\|_2 &= \left\| \frac{1}{n} (\mathbf{A}^{\mathbf{X}_k} \otimes \mathbf{I}_M)^\top \mathbf{v} \right\|_2 \\ &= \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_i^{\mathbf{X}_k} \otimes \mathbf{I}_M) \mathbf{v}_i \right\|_2 = \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{a}_i^{\mathbf{X}_k} \mathbf{v}_i^\top \right\|_{\mathbb{F}}, \end{aligned} \quad (75)$$

and by Lemma B.6, we further have

$$\left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{v} \right\|_2 \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{a}_i^{\mathbf{X}_k}\|_2^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i\|_2^2}. \quad (76)$$

By Lemma B.14 with

$$\begin{aligned} \Gamma_g(j) &= \sigma_{\max, 0} \mathcal{Q}_{n, p, \delta}^2 \text{ for all } j \in [p] \\ \Gamma_e &= \sigma_{jr} \mathcal{Q}_{n, \delta}^1 \\ \Gamma_\phi(m) &= d_{jm} \sqrt{\frac{(1/c_1) \log(c_2/\delta)}{n}} \end{aligned}$$

because of  $\mathcal{A}_1$ ,  $\mathcal{A}_2$  and Assumption 4.4, we have

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i\|^2 \leq \Gamma_v, \quad (77)$$

where  $\Gamma_v$  is defined by (68). And combine with  $\mathcal{A}_4$ , we have

$$\left\| \frac{2}{n} \mathbf{Z}_k^\top \mathbf{v} \right\|_2 \leq 2 \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{a}_i^{\mathbf{X}_k}\|_2^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i\|_2^2} \leq 2 \sqrt{\text{tr}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})} \sqrt{\mathcal{Q}_{n,p-1,\delta}^2} \cdot \Gamma_v^{1/2}.$$

Finally, to bound the first term in (74), by similar arguments as (75) and (76), we have

$$\left\| \frac{1}{n} (\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}} \right\|_2 \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{a}}_i^{\mathbf{X}_k} - \mathbf{a}_i^{\mathbf{X}_k}\|_2^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{u}}_i\|_2^2}. \quad (78)$$

Since we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{a}}_i^{\mathbf{X}_k} - \mathbf{a}_i^{\mathbf{X}_k}\|_2^2 &\leq \left( \frac{1}{n} \sum_{i=1}^n \|g_{ik}\|^2 \right) \sum_{m=1}^M \|\hat{\phi}_m - \phi_m\|^2 \\ &\leq \sigma_{\max,0} \mathcal{Q}_{n,p,\delta}^2 \frac{(1/c_1) \log(c_2/\delta)}{n} d_{js}^2(M) \end{aligned}$$

by  $\mathcal{A}_1$  and Assumption 4.4, and

$$\begin{aligned} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{u}}_i\|_2^2} &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{r}_i + \mathbf{w}_i\|_2^2} + \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i\|_2^2} \\ &\leq \sqrt{\Xi_4(M)} \sqrt{\mathcal{Q}_{n,\delta}^1} + \Gamma_v^{1/2} \end{aligned}$$

by  $\mathcal{A}_3$  and (77), then by (78), we have

$$\begin{aligned} \left\| \frac{1}{n} (\hat{\mathbf{Z}}_k - \mathbf{Z}_k)^\top \tilde{\mathbf{u}} \right\|_2 &\leq \sqrt{\sigma_{\max,0}} \sqrt{\mathcal{Q}_{n,p,\delta}^2} \sqrt{\frac{(1/c_1) \log(c_2/\delta)}{n}} \\ &\quad d_{js}(M) \left( \sqrt{\Xi_4(M)} \sqrt{\mathcal{Q}_{n,\delta}^1} + \Gamma_v^{1/2} \right). \end{aligned}$$

Thus, combine the bounds of three terms in (74) and by the definition of  $\check{\lambda}(n, p, M, \delta)$  in (69), we have

$$\frac{2}{n} \max_{k \in [p-1]} \|\mathbf{Z}_k^\top \mathbf{v}\|_2 \leq \check{\lambda}(n, p, M, \delta).$$

We then prove that  $\cap_{i=1}^7 \mathcal{A}_i$  implies that  $\|\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty \leq \check{\Gamma}(n, p, M, \delta)$ . Note that

$$\|\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty \leq \|\hat{\boldsymbol{\Sigma}}^{\mathbf{X}} - \hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}}\|_\infty + \|\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}\|_\infty,$$

and by  $\mathcal{A}_6$ , we thus only need to prove that  $\left\| \tilde{\Sigma}^{\mathbf{X}} - \hat{\Sigma}_n^{\mathbf{X}} \right\|_{\infty} \leq \check{\Gamma}(n, p, M, \delta) - \Gamma(n, p, M, \delta)$ . Note that

$$\begin{aligned} \left\| \tilde{\Sigma}^{\mathbf{X}} - \hat{\Sigma}_n^{\mathbf{X}} \right\|_{\infty} &= \max_{1 \leq k, k' \leq p-1} \left\| \frac{1}{n} (\hat{\mathbf{A}}^{\mathbf{X}_k})^{\top} \hat{\mathbf{A}}^{\mathbf{X}_{k'}} - \frac{1}{n} (\mathbf{A}^{\mathbf{X}_k})^{\top} \mathbf{A}^{\mathbf{X}_{k'}} \right\|_{\infty} \\ &= \max_{1 \leq k, k' \leq p-1} \max_{1 \leq m, m' \leq M} \left| \frac{1}{n} \sum_{i=1}^n \hat{a}_{im}^{X_k} \hat{a}_{im'}^{X_{k'}} - a_{im}^{X_k} a_{im'}^{X_{k'}} \right|, \end{aligned} \quad (79)$$

and

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n \hat{a}_{im}^{X_k} \hat{a}_{im'}^{X_{k'}} - a_{im}^{X_k} a_{im'}^{X_{k'}} \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k} - a_{im}^{X_k}) (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}}) + \frac{1}{n} \sum_{i=1}^n a_{im}^{X_k} (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}}) \right. \\ & \quad \left. + \frac{1}{n} \sum_{i=1}^n a_{im'}^{X_{k'}} (\hat{a}_{im}^{X_k} - a_{im}^{X_k}) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k} - a_{im}^{X_k}) (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}}) \right| + \left| \frac{1}{n} \sum_{i=1}^n a_{im}^{X_k} (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}}) \right| \\ & \quad + \left| \frac{1}{n} \sum_{i=1}^n a_{im'}^{X_{k'}} (\hat{a}_{im}^{X_k} - a_{im}^{X_k}) \right| \\ &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k} - a_{im}^{X_k})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}})^2} \\ & \quad + \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im'}^{X_{k'}} - a_{im'}^{X_{k'}})^2} \\ & \quad + \sqrt{\frac{1}{n} \sum_{i=1}^n (a_{im'}^{X_{k'}})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k} - a_{im}^{X_k})^2}. \end{aligned} \quad (80)$$

By Lemma B.12 with

$$\begin{aligned} \Gamma_g(j) &= \sigma_{\max, 0} \mathcal{Q}_{n, p, \delta}^2 \text{ for all } j \in [p], \\ \Gamma_{\phi}(m) &= d_{jm} \sqrt{\frac{(1/c_1) \log(c_2/\delta)}{n}}, \end{aligned}$$

we have

$$\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im}^{X_k} - a_{im}^{X_k})^2 \leq d_{j, \max}^2 \sigma_{\max, 0} \mathcal{Q}_{n, p, \delta}^2 \frac{(1/c_1) \log(c_2/\delta)}{n}$$

for all  $m \in [M]$  and  $k \in [p-1]$ . And by  $\mathcal{A}_7$ , we have

$$\frac{1}{n} \sum_{i=1}^n \|a_{im}^{X_k}\|_2^2 \leq \Xi_2(M) \mathcal{Q}_{n,p,\delta}^3$$

for all  $m \in [M]$  and  $k \in [p-1]$ . Combine with (79) and (80), we have

$$\begin{aligned} & \left\| \tilde{\Sigma}^{\mathbf{X}} - \hat{\Sigma}_n^{\mathbf{X}} \right\|_{\infty} \\ & \leq d_{j,\max}^2 \sigma_{\max,0} \mathcal{Q}_{n,p,\delta}^2 \frac{(1/c_1) \log(c_2/\delta)}{n} \\ & \quad + 2\sqrt{\Xi_2(M)} d_{j,\max} \sqrt{\sigma_{\max,0}} \sqrt{\mathcal{Q}_{n,p,\delta}^2} \sqrt{\frac{(1/c_1) \log(c_2/\delta)}{n}} \cdot \sqrt{\mathcal{Q}_{n,p,M\delta}^3} \\ & = \check{\Gamma}(n, p, M, \delta) - \Gamma(n, p, M, \delta) \leq \check{\Gamma}(n, p, M, \delta). \end{aligned}$$

Therefore, by (73) and Lemma B.3, Lemma B.5, Lemma B.9 and Lemma B.10, we have

$$\begin{aligned} & \mathbb{P} \left\{ \frac{2}{n} \max_{k \in [p-1]} \|(\mathbf{Z}^{\top} \mathbf{u})_k\|_2 \leq \check{\lambda}(n, p, M, \delta) \text{ and } \left\| \tilde{\Sigma}_n^{\mathbf{X}} - \Sigma^{\mathbf{X}} \right\|_{\infty} \leq \check{\Gamma}(n, p, M, \delta) \right\} \\ & \geq 1 - \mathbb{P} \left\{ \cup_{i=1}^7 \bar{\mathcal{A}}_i \right\} \\ & \geq 1 - 8\delta. \end{aligned}$$

Finally, we only need  $n$  large enough such that  $\check{\Gamma}(n, p, M, \delta) \leq \kappa(M)/(32M^2s)$  and  $\check{\chi}(n, p, M, \delta) \leq \nu(M)/3$  where  $\check{\chi}(n, p, M, \delta)$  is defined in (38). After dropping  $\log(n)$  term, to satisfy the first condition, we need

$$n \geq O \left( \frac{M^4 s^2 \log(p^2 M^2 / \delta)}{\kappa^2(M)} \right),$$

and to satisfy the second condition, we need

$$n \geq O \left( \frac{1}{\kappa(M) \nu^2(M)} \max \left\{ sM^2, s \log(p/\delta), sM^2 \omega^2(M) \log(M^2 p/\delta), \right. \right. \\ \left. \left. s^3 M^2 (\log(1/\delta)^2), s(d_{j_0}^2 - d_{j_s}^2(M) \Phi^2(M)) \right\} \right).$$

Then combine the above results and note that decreasing  $8\delta$  to  $\delta$  doesn't affect the asymptotic order of  $n$ , we have the final result.  $\square$

## Appendix B: Useful lemmas

Recall that  $\mathcal{C}_{n,\delta}$  is defined in (34) and  $\mathcal{Q}_{n,\delta}^1$ ,  $\mathcal{Q}_{n,p,\delta}^2$ , and  $\mathcal{Q}_{n,p,M,\delta}^1$  are defined in (67).

**Lemma B.1.** *Let  $\delta \in (0, 1]$ . Then  $\max_{k \in [p-1]} \|\mathbf{A}^{X_k}\|_2 / \sqrt{n} \leq \mathcal{C}_{n,\delta}$  with probability at least  $1 - \delta$ .*

*Proof.* The result follows directly from Theorem 6.1 of [63] and a union bound.  $\square$

**Lemma B.2.** *Suppose  $\max_{k \in [p-1]} \|\mathbf{A}^{\mathbf{X}_k}\|_2 / \sqrt{n} \leq \mathcal{C}_{n, \delta_1}$  for some  $0 < \delta_1 \leq 1$ . Then there exists a constant  $c > 0$  such that*

$$\begin{aligned} \max_{k \in [p-1]} \left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{u} \right\|_2 &\leq \frac{M \mathcal{C}_{n, \delta_1} \sqrt{\Xi_1(M)}}{\sqrt{n}} + 2 \mathcal{C}_{n, \delta_1} \sqrt{\Xi_1(M)} \sqrt{\frac{\log(4(p-1)/\delta_2)}{n}} \\ &+ \omega(M) \left\{ 1 + 7\sqrt{3} \frac{M \sqrt{\log(6(p-1)/\delta_2) + 2 \log M}}{\sqrt{n}} + \right. \\ &\quad \left. \frac{8c M (\log(2n)) (\log(6(p-1)/\delta_2) + 2 \log M)}{3n} \right\} \end{aligned}$$

holds with probability at least  $1 - \delta_2$ ,  $\delta_2 \in (0, 1]$ .

*Proof.* Note that  $\|\mathbf{Z}_k\|_2 = \|\mathbf{A}^{\mathbf{X}_k} \otimes \mathbf{I}_M\|_2 = \|\mathbf{A}^{\mathbf{X}_k}\|_2$ . Thus  $\|\mathbf{Z}_k / \sqrt{n}\|_2 \leq \mathcal{C}_{n, \delta_1}$ ,  $k \in [p-1]$ , since  $\max_{k \in [p-1]} \|\mathbf{A}^{\mathbf{X}_k}\|_2 / \sqrt{n} \leq \mathcal{C}_{n, \delta_1}$ . Recall that  $\mathbf{u} = \mathbf{w} + \mathbf{r}$ , where  $\mathbf{u}$  is defined in (57),  $\mathbf{w} = (\mathbf{w}_1^\top, \mathbf{w}_2^\top, \dots, \mathbf{w}_n^\top)^\top \in \mathbb{R}^{nM}$  and  $\mathbf{r} = (\mathbf{r}_1^\top, \mathbf{r}_2^\top, \dots, \mathbf{r}_n^\top)^\top \in \mathbb{R}^{nM}$ . Both  $\mathbf{w}_i$  and  $\mathbf{r}_i$  are Gaussian vectors with mean zero, and covariance matrices  $\Sigma^{\mathbf{w}}$  and  $\Sigma^{\mathbf{r}}$  respectively for all  $i \in [n]$ , where we dropped the superscript  $M$  to simplify the notation, and  $\mathbf{w}_i \perp \mathbf{r}_i$ .

Fix  $k \in [p-1]$ . Let  $\mathbf{r}_i = \mathbf{r}_i^1 + \mathbf{r}_i^2$ , where

$$\mathbf{r}_i^1 = \mathbf{r}_i - \Sigma^{r, \mathbf{X}_k} (\Sigma_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k} \quad \text{and} \quad \mathbf{r}_i^2 = \Sigma^{r, \mathbf{X}_k} (\Sigma_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k}.$$

Since both  $\mathbf{r}_i$  and  $\mathbf{a}_i^{\mathbf{X}_k}$  are Gaussian vectors and  $\text{Cov}(\mathbf{a}_i^{\mathbf{X}_k}, \mathbf{r}_i^1) = \mathbf{0}$ , we have  $\mathbf{a}_i^{\mathbf{X}_k} \perp \mathbf{r}_i^1$ . Then

$$\left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{u} \right\|_2 \leq \left\| \frac{1}{n} \mathbf{Z}_k^\top (\mathbf{w} + \mathbf{r}^1) \right\|_2 + \left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{r}^2 \right\|_2. \quad (81)$$

We upper bound the first term on the right hand side of (81). Let  $\boldsymbol{\xi}_i = \mathbf{w}_i + \mathbf{r}_i^1$ . Then  $\boldsymbol{\xi}_i$  is a Gaussian vector with mean zero and covariance matrix  $\Sigma_{kk}^{\boldsymbol{\xi}}$ , with  $\Sigma_{kk}^{\boldsymbol{\xi}} = \Sigma^{\mathbf{w}} + \Sigma^{\mathbf{r}} - \Sigma^{r, \mathbf{X}_k} (\Sigma_{kk}^{\mathbf{X}})^{-1} \Sigma^{\mathbf{X}_k, r}$ . Note that  $\boldsymbol{\xi} \perp \mathbf{Z}_k$  and we establish a bound conditional on  $\mathbf{Z}_k$  first.

For any pair  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^{nM}$ , we have

$$\begin{aligned} \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\theta}}{n} \right\|_2 - \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\theta}'}{n} \right\|_2 &\leq \frac{1}{n} \|\mathbf{Z}_k^\top (\boldsymbol{\theta} - \boldsymbol{\theta}')\|_2 \\ &\leq \frac{1}{n} \|\mathbf{Z}_k\|_2 \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 \leq \frac{\mathcal{C}_{n, \delta_1}}{\sqrt{n}} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2. \end{aligned}$$

Therefore,  $\boldsymbol{\theta} \mapsto \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\theta}}{n} \right\|_2$  is a Lipschitz function with a Lipschitz constant  $\frac{\mathcal{C}_{n, \delta_1}}{\sqrt{n}}$ . Note that  $\boldsymbol{\xi}$  has a strongly log-concave distribution with the parameter

$$(\rho_{\max}(\mathbf{I}_M \otimes \Sigma_{kk}^{\boldsymbol{\xi}}))^{-1} = (\rho_{\max}(\Sigma_{kk}^{\boldsymbol{\xi}}))^{-1}.$$

Then, by Theorem 3.16 of [63], we have

$$\mathbb{P} \left\{ \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\xi}}{n} \right\|_2 \geq \mathbb{E} \left[ \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\xi}}{n} \right\|_2 \right] + \Delta \right\} \leq 2 \exp \left( - \frac{n\Delta^2}{4C_{n,\delta_1}^2 \rho_{\max}(\boldsymbol{\Sigma}_k^\xi)} \right) \quad (82)$$

for all  $\Delta > 0$ .

Next, we bound  $\mathbb{E} \left[ \frac{1}{n} \left\| \mathbf{Z}_k^\top \boldsymbol{\xi} \right\|_2 \right]$  using the Sudakov-Fernique inequality [Theorem 2.2.3 of 1]. For any vector  $\boldsymbol{\theta} \in \mathbb{R}^{M^2}$ , let  $\zeta_{\boldsymbol{\theta}} = n^{-1} \langle \boldsymbol{\theta}, \mathbf{Z}_k^\top \boldsymbol{\xi} \rangle$ . Then  $\frac{1}{n} \left\| \mathbf{Z}_k^\top \boldsymbol{\xi} \right\|_2 = \max_{\|\boldsymbol{\theta}\|_2=1} \zeta_{\boldsymbol{\theta}}$ . For any two vectors  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^{M^2}$ , we have

$$\begin{aligned} \mathbb{E} [(\zeta_{\boldsymbol{\theta}} - \zeta_{\boldsymbol{\theta}'})^2] &= \mathbb{E} \left[ \left( \frac{1}{n} \langle \boldsymbol{\theta} - \boldsymbol{\theta}', \mathbf{Z}_k^\top \boldsymbol{\xi} \rangle \right)^2 \right] \\ &= \frac{1}{n^2} \mathbb{E} [(\boldsymbol{\theta} - \boldsymbol{\theta}')^\top \mathbf{Z}_k^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top \mathbf{Z}_k (\boldsymbol{\theta} - \boldsymbol{\theta}')] \\ &\leq \frac{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}{n} \frac{\|\mathbf{Z}_k\|_2^2}{n} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2 \leq \frac{C_{n,\delta_1}^2 \rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}{n} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2. \end{aligned}$$

We define another Gaussian process,

$$\tilde{\zeta}_{\boldsymbol{\theta}} = \left( C_{n,\delta_1} \sqrt{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)/n} \right) \cdot \langle \boldsymbol{\theta}, \boldsymbol{\epsilon} \rangle,$$

where  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I}_{M^2})$ . For any pair  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^{M^2}$ , we have

$$\mathbb{E} [(\tilde{\zeta}_{\boldsymbol{\theta}} - \tilde{\zeta}_{\boldsymbol{\theta}'})^2] = \frac{C_{n,\delta_1}^2 \rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}{n} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2 \geq \mathbb{E} [(\zeta_{\boldsymbol{\theta}} - \zeta_{\boldsymbol{\theta}'})^2].$$

Then, by the Sudakov-Fernique inequality, we have

$$\mathbb{E} \left[ \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\xi}}{n} \right\|_2 \right] = \mathbb{E} \left[ \max_{\|\boldsymbol{\theta}\|_2=1} \zeta_{\boldsymbol{\theta}} \right] \leq \mathbb{E} \left[ \max_{\|\boldsymbol{\theta}\|_2=1} \tilde{\zeta}_{\boldsymbol{\theta}} \right].$$

On the other hand,

$$\begin{aligned} \mathbb{E} \left[ \max_{\|\boldsymbol{\theta}\|_2=1} \tilde{\zeta}_{\boldsymbol{\theta}} \right] &= \frac{C_{n,\delta_1} \sqrt{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}}{\sqrt{n}} \mathbb{E} [\|\boldsymbol{\epsilon}\|_2] \leq \frac{C_{n,\delta_1} \sqrt{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}}{\sqrt{n}} \sqrt{\mathbb{E} [\|\boldsymbol{\epsilon}\|_2^2]} \\ &= \frac{MC_{n,\delta_1} \sqrt{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}}{\sqrt{n}}. \end{aligned}$$

Combining (82) with the above upper bound, we have

$$\mathbb{P} \left\{ \left\| \frac{\mathbf{Z}_k^\top \boldsymbol{\xi}}{n} \right\|_2 \geq \frac{MC_{n,\delta_1} \sqrt{\rho_{\max}(\boldsymbol{\Sigma}_k^\xi)}}{\sqrt{n}} + \Delta \right\}$$

$$\leq 2 \exp\left(-\frac{n\Delta^2}{4C_{n,\delta_1}^2 \rho_{\max}(\boldsymbol{\Sigma}_{\mathbf{k}}^{\xi})}\right), \quad \Delta > 0.$$

The above inequality holds for any  $\mathbf{Z}_{\mathbf{k}}$  such that  $\|\mathbf{Z}_{\mathbf{k}}/\sqrt{n}\|_2 \leq C_{n,\delta_1}$  and, therefore, is valid unconditionally as well. Since  $\Xi_1(M) = \max_{k \in [p-1]} \rho_{\max}(\boldsymbol{\Sigma}_{\mathbf{k}}^{\xi})$ , it follows from the union bound that

$$\begin{aligned} \mathbb{P}\left\{\max_{k \in [p-1]} \left\|\frac{\mathbf{Z}_{\mathbf{k}}^{\top} \boldsymbol{\xi}}{n}\right\|_2 \geq \frac{MC_{n,\delta_1} \sqrt{\Xi_1(M)}}{\sqrt{n}} + \Delta\right\} \\ \leq 2(p-1) \exp\left(-\frac{n\Delta^2}{4C_{n,\delta_1}^2 \Xi_1(M)}\right), \quad \Delta > 0. \end{aligned} \quad (83)$$

Next, we upper bound the second term in (81). For any  $k \in [p-1]$ , we have

$$\begin{aligned} & \left\|\frac{1}{n} \mathbf{Z}_{\mathbf{k}}^{\top} \mathbf{r}^2\right\|_2 \\ &= \left\|\frac{1}{n} \sum_{i=1}^n (\mathbf{a}_i^{\mathbf{X}_k} \otimes \mathbf{I}_M) \mathbf{r}_i^2\right\|_2 = \left\|\frac{1}{n} \sum_{i=1}^n \mathbf{r}_i^2 (\mathbf{a}_i^{\mathbf{X}_k})^{\top}\right\|_{\mathbb{F}} \\ &= \left\|\frac{1}{n} \sum_{i=1}^n \boldsymbol{\Sigma}^{r, \mathbf{X}_k} (\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k} (\mathbf{a}_i^{\mathbf{X}_k})^{\top}\right\|_{\mathbb{F}} \\ &\leq \left\|\frac{1}{n} \sum_{i=1}^n \boldsymbol{\Sigma}^{r, \mathbf{X}_k} (\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k} (\mathbf{a}_i^{\mathbf{X}_k})^{\top} - \boldsymbol{\Sigma}^{r, \mathbf{X}_k}\right\|_{\mathbb{F}} + \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \\ &= \left\|\boldsymbol{\Sigma}^{r, \mathbf{X}_k} \left\{\frac{1}{n} \sum_{i=1}^n (\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k} (\mathbf{a}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\}\right\|_{\mathbb{F}} + \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \\ &\leq \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \left\|\frac{1}{n} \sum_{i=1}^n (\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1} \mathbf{a}_i^{\mathbf{X}_k} (\mathbf{a}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\|_{\mathbb{F}} + \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \\ &= \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \left\|\frac{1}{n} \sum_{i=1}^n \left((\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1/2} \mathbf{a}_i^{\mathbf{X}_k}\right) \left(\left((\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1/2} \mathbf{a}_i^{\mathbf{X}_k}\right)\right)^{\top} - \mathbf{I}_M\right\|_{\mathbb{F}} \\ &\quad + \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \\ &\triangleq \|\boldsymbol{\Sigma}^{r, \mathbf{X}_k}\|_{\mathbb{F}} \left\{\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^{\mathbf{X}_k} (\mathbf{b}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\|_{\mathbb{F}} + 1\right\}, \end{aligned} \quad (84)$$

where  $\mathbf{b}_i^{\mathbf{X}_k} = (\boldsymbol{\Sigma}_{kk}^{\mathbf{X}})^{-1/2} \mathbf{a}_i^{\mathbf{X}_k} \sim N(0, \mathbf{I}_M)$ . Since

$$\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^{\mathbf{X}_k} (\mathbf{b}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\|_{\mathbb{F}} \leq M \left\|\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^{\mathbf{X}_k} (\mathbf{b}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\|_{\infty}\right\|, \quad (85)$$

we only need to bound  $\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^{\mathbf{X}_k} (\mathbf{b}_i^{\mathbf{X}_k})^{\top} - \mathbf{I}_M\right\|_{\infty}$ . We use Theorem 4.1 in [31] to bound this term. We first check the conditions therein. Note that for

$b_{im}^{X_k}$  is a standard normal random variable for all  $i \in [n]$  and  $m \in [M]$ , we have  $(b_{im}^{X_k})^2$  to be chi-squared distributed with degree of freedom 1. By the moment generating function of chi-square distribution, we have

$$\mathbb{E} \left[ \exp\{\eta(b_{im}^{X_k})^2\} \right] = \frac{1}{\sqrt{1-2\eta}} \quad \text{for all } 0 < \eta < 1/2.$$

Thus, we have

$$\begin{aligned} \|b_{im}^{X_k}\|_{\psi_2} &= \inf\{\eta > 0 : \mathbb{E} \left[ \psi_2(|b_{im}^{X_k}|/\eta) \right] \leq 1\} \\ &= \inf\{\eta > 0 : \mathbb{E} \left[ \exp((b_{im}^{X_k})^2/\eta^2) \right] \leq 2\} \\ &= \inf\{\eta > 0 : 1/\sqrt{1-2/\eta^2} \leq 2\} \\ &= \frac{2\sqrt{2}}{\sqrt{3}}. \end{aligned}$$

Furthermore, we have

$$\|b_i^{X_k}\|_{M, \psi_2} = \max_{m \in [M]} \|b_{im}^{X_k}\|_{\psi_2} = \frac{2\sqrt{2}}{\sqrt{3}}.$$

for all  $i \in [n]$ . In addition, note that

$$\text{Var} \left( b_{im}^{X_k} b_{im'}^{X_k} \right) \leq \mathbb{E} \left[ (b_{im}^{X_k})^2 (b_{im'}^{X_k})^2 \right] \leq \left( \mathbb{E} \left[ (b_{im}^{X_k})^4 \right] \right)^{1/2} \left( \mathbb{E} \left[ (b_{im'}^{X_k})^4 \right] \right)^{1/2} = 3$$

for all  $m, m' \in [M]$ , thus we have

$$\max_{m, m' \in [M]} \left\{ \frac{1}{n} \sum_{i=1}^n \text{Var} \left( b_{im}^{X_k} b_{im'}^{X_k} \right) \right\} \leq 3.$$

Therefore, by Theorem 4.1 in [31], we have

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^{X_k} \left( \mathbf{b}_i^{X_k} \right)^\top - \mathbf{I}_M \right\|_\infty > 7\sqrt{3} \sqrt{\frac{\Delta + 2 \log M}{n}} + \frac{8c (\log(2n))(\Delta + 2 \log M)}{3n} \right\} \leq 3e^{-\Delta},$$

for all  $\Delta > 0$ , where  $c$  is a constant. Thus, combining the above equation with (84) and (85), we have

$$\begin{aligned} &\mathbb{P} \left\{ \left\| \frac{1}{n} \mathbf{Z}_k^\top \mathbf{r}^2 \right\|_2 > \right. \\ &\left. \left\| \Sigma^{\mathbf{r}, \mathbf{X}_k} \right\|_F \left( 1 + 7\sqrt{3} \frac{M \sqrt{\Delta + 2 \log M}}{\sqrt{n}} + \frac{8c M (\log(2n))(\Delta + 2 \log M)}{3n} \right) \right\} \\ &\leq 3e^{-\Delta}. \end{aligned}$$



Let  $\omega(M) = \max_{k \in [p-1]} \|\boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{X}_k}\|_{\mathbb{F}}$ , then

$$\begin{aligned} \mathbb{P} \left\{ \max_{k \in [p-1]} \left\| \frac{1}{n} \mathbf{Z}_k^{\top} \mathbf{r}^2 \right\|_2 > \omega(M) \left( 1 + 7\sqrt{3} \frac{M\sqrt{\Delta + 2\log M}}{\sqrt{n}} \right. \right. \\ \left. \left. + \frac{8c}{3} \frac{M(\log(2n))(\Delta + 2\log M)}{n} \right) \right\} \\ \leq 3(p-1)e^{-\Delta}. \end{aligned} \quad (86)$$

The result follows by combining (81), (83), and (86).  $\square$

**Lemma B.3.** Let  $\lambda_n = \tilde{\lambda}(n, p, M, \delta)$ , where

$$\begin{aligned} \tilde{\lambda}(n, p, M, \delta) = 2\mathcal{C}_{n, \delta} \left( \frac{M\sqrt{\Xi_1(M)}}{\sqrt{n}} + 2\sqrt{\Xi_1(M)} \sqrt{\frac{\log(4(p-1)/\delta)}{n}} \right) \\ + 2\omega(M) \left\{ 7\sqrt{3} \frac{M\sqrt{\log(6(p-1)/\delta) + 2\log M}}{\sqrt{n}} \right. \\ \left. + \frac{8Mc(\log(2n))(\log(6(p-1)/\delta) + 2\log M)}{3n} \right\}, \end{aligned}$$

is defined in (54), then (62) holds with probability at least  $1 - 2\delta$ . That is, we have

$$\lambda_n \geq \frac{2}{n} \max_{k \in [p-1]} \|(\mathbf{Z}^{\top} \mathbf{u})_k\|_2$$

hold with probability at least  $1 - 2\delta$ .

*Proof.* The result follows directly from Lemma B.1, Lemma B.2 and Bonferroni inequality.  $\square$

**Lemma B.4.**

$$\mathbb{P} \left\{ \sqrt{\sum_{k=1}^{p-1} \left\| \hat{\mathbf{B}}_k - \mathbf{B}_k^* \right\|_F^2} \leq \chi(n, p, M, \delta) + 12\sqrt{s} \frac{\omega(M)}{\sqrt{\kappa(M)}} \right\} \geq 1 - 3\delta,$$

where  $\kappa(M)$  is defined in (28) and  $\chi(n, p, M, \delta)$  is defined in (32).

*Proof.* By Lemma B.3 and (64), we can get an error bound for  $\|\hat{\beta} - \beta^*\|_2$  by applying Theorem 1 in [47]. Note that by Lemma C.4 of [68], we have the subspace compatibility constant defined in Definition 3 of [47] to be  $\sqrt{s}$ . Then, for all  $0 < \delta \leq 1$ , when  $n$  is large enough such that  $\Gamma(n, p, M, \delta) \leq \kappa(M)/(32M^2s)$ , where  $\Gamma(n, p, M, \delta)$  is defined in (88), and let  $\lambda_n = \tilde{\lambda}(n, p, M, \delta)$ , where  $\tilde{\lambda}(n, p, M, \delta)$  is defined in (33), then by Corollary 1 in [47], Lemma B.3, (64), and a union bound, we have the desired result.  $\square$

**Lemma B.5.** There exists a constant  $c$  such that

$$\begin{aligned} \left\| \widehat{\Sigma}_n^{\mathbf{X}} - \Sigma^{\mathbf{X}} \right\|_{\infty} &\leq 7\sqrt{3}K_0 \sqrt{\frac{\log(3/\delta) + 2\log((p-1)M)}{n}} \\ &\quad + \frac{8cK_0 \log(2n)(\log(3/\delta) + 2\log((p-1)M))}{3n} \end{aligned} \quad (87)$$

holds with probability at least  $1 - \delta$  where  $K_0$  is defined by (27).

*Proof.* The result follows Theorem 4.1 in [31]. In order to apply the theorem, we need to check the conditions therein. First, we bound  $\max_{k \in [p-1], m \in [M]} \|a_{im}^{X_k}\|_{\psi_2}$ . Note that for  $\zeta \sim N(0, \sigma^2)$ , we have  $(\zeta/\sigma)^2$  is chi-square distributed with degree of freedom 1. By the moment generating function of chi-square distribution, we have

$$\mathbb{E} [\exp\{\eta(\zeta/\sigma)^2\}] = \frac{1}{\sqrt{1-2\eta}} \quad \text{for all } 0 < \eta < 1/2.$$

Thus, we have

$$\mathbb{E} [\exp(\zeta^2/t^2)] = \mathbb{E} [\exp\{(\sigma^2/t^2) \cdot (\zeta^2/\sigma^2)\}] = \frac{1}{\sqrt{1-2\sigma^2/t^2}}$$

for all  $t > \sqrt{2}\sigma$ . Let  $1/\sqrt{1-2\sigma^2/t^2} \leq 2$ , we have  $t \geq (2\sqrt{2})/(\sqrt{3}\sigma)$ . Thus, we have

$$\begin{aligned} \|\zeta\|_{\psi_2} &= \inf\{t > 0 : \mathbb{E}[\psi_2(|\zeta|/t)] \leq 1\} \\ &= \inf\{t > 0 : \mathbb{E}[\exp(\zeta^2/t^2)] \leq 2\} \\ &= \inf\{t > 0 : 1/\sqrt{1-2\sigma^2/t^2} \leq 2\} \\ &= \frac{2\sqrt{2}}{\sqrt{3}}\sigma. \end{aligned}$$

Based on the above result, we have for any  $i \in [n]$ ,

$$\begin{aligned} \|\mathbf{a}_i^{\mathbf{X}_k}\|_{M, \psi_2} &= \max_{k \in [p-1], m \in [M]} \|a_{im}^{X_k}\|_{\psi_2} \\ &= \frac{2\sqrt{2}}{\sqrt{3}} \max_{k \in [p-1], m \in [M]} \sqrt{\Sigma_{kk,mm}^X} \\ &= \frac{2\sqrt{2}}{\sqrt{3}} \sqrt{K_0}, \end{aligned}$$

where  $K_0$  is defined by (27). We then bound  $\max_{k, k' \in [p], m, m' \in [M]} \text{Var}(a_{im}^{X_k} a_{im'}^{X_{k'}})$ . This followed by

$$\begin{aligned} \text{Var}(a_{im}^{X_k} a_{im'}^{X_{k'}}) &\leq \mathbb{E}[(a_{im}^{X_k})^2 (a_{im'}^{X_{k'}})^2] \leq \left(\mathbb{E}[(a_{im}^{X_k})^4]\right)^{1/2} \left(\mathbb{E}[(a_{im'}^{X_{k'}})^2]\right)^{1/2} \\ &= (9(\Sigma_{kk,mm}^X)^2 (\Sigma_{k'k',m'm'}^X)^2)^{1/2} \\ &\leq 3K_0^2. \end{aligned}$$

Thus, the final result is derived by applying Theorem 4.1 in [31] with  $K_{n,p} = (2\sqrt{2}/\sqrt{3})\sqrt{K_0}$  and  $A_{n,p}^2 = 3K_0^2$ .  $\square$

**Lemma B.6.** For  $\mathbf{x}_i, \mathbf{y}_i \in \mathbb{R}^M$ ,  $i \in [n]$ , we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i^\top \right\|_F \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|_2^2} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_i\|_2^2}.$$

*Proof.* For any  $m, m' \in [M]$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n x_{im} y_{im'} \right|^2 \leq \left( \frac{1}{n} \sum_{i=1}^n x_{im}^2 \right) \left( \frac{1}{n} \sum_{i=1}^n y_{im'}^2 \right).$$

Therefore,

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i^\top \right\|_F^2 &= \sum_{m, m'=1}^M \left| \frac{1}{n} \sum_{i=1}^n x_{im} y_{im'} \right|^2 \\ &\leq \sum_{m, m'=1}^M \left( \frac{1}{n} \sum_{i=1}^n x_{im}^2 \right) \left( \frac{1}{n} \sum_{i=1}^n y_{im'}^2 \right) \\ &= \left( \frac{1}{n} \sum_{i=1}^n \sum_{m=1}^M x_{im}^2 \right) \left( \frac{1}{n} \sum_{i=1}^n \sum_{m'=1}^M y_{im'}^2 \right) \\ &= \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 \right) \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_i\|_2^2 \right), \end{aligned}$$

and the result immediately follows.  $\square$

**Lemma B.7.** Let

$$\Gamma(n, p, M, \delta) = 7\sqrt{3}K_0 \sqrt{\frac{\log(3/\delta) + 2 \log((p-1)M)}{n}} + \frac{8cK_0 \log(2n)(\log(3/\delta) + 2 \log((p-1)M))}{3n}, \quad (88)$$

then when  $\Gamma(n, p, M, \delta) \leq \kappa(M)/(32M^2s)$ , we have

$$\mathbb{P} \left\{ \frac{1}{2n} \|\mathbf{Z} \Delta \boldsymbol{\beta}\|_2^2 \geq \frac{\kappa}{4} \|\Delta \boldsymbol{\beta}\|_2^2 \text{ for all } \Delta \boldsymbol{\beta} \in \mathbb{C}(\mathcal{N}_j) \right\} \geq 1 - \delta,$$

where  $\mathbb{C}(\mathcal{N}_j)$  is defined in (63).

*Proof.* We want to first prove

$$\frac{1}{2n} \|\mathbf{Z} \Delta \boldsymbol{\beta}\|_2^2 \geq \kappa_{\mathcal{L}} \|\Delta \boldsymbol{\beta}\|_2^2 - \tau_{\mathcal{L}}(\boldsymbol{\beta}^*) \quad \text{for all } \Delta \boldsymbol{\beta} \in \mathbb{C}(\mathcal{N}_j),$$

where  $\kappa_{\mathcal{L}} > 0$  is a positive constant,  $\tau_{\mathcal{L}}(\boldsymbol{\beta}^*) > 0$ , and  $\mathbb{C}(\mathcal{N}_j)$  is defined in (63).

Since  $\mathbf{Z} = \mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M$ , thus for any  $\boldsymbol{\theta} \in \mathbb{R}^{(p-1)M^2}$ , we have

$$\begin{aligned} \frac{1}{n} \|\mathbf{Z}\boldsymbol{\theta}\|_2^2 &= \frac{1}{n} \|(\mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M)\boldsymbol{\theta}\|_2^2 \\ &= \frac{1}{n} \boldsymbol{\theta}^\top (\mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M)^\top (\mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta} \\ &= \frac{1}{n} \boldsymbol{\theta}^\top ((\mathbf{A}^{\mathbf{X}})^\top \otimes \mathbf{I}_M) (\mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta} \\ &= \frac{1}{n} \boldsymbol{\theta}^\top ((\mathbf{A}^{\mathbf{X}})^\top \mathbf{A}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta} \\ &= \boldsymbol{\theta}^\top (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta}, \end{aligned}$$

where  $\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} = \frac{1}{n} (\mathbf{A}^{\mathbf{X}})^\top \mathbf{A}^{\mathbf{X}}$ . We then further have

$$\begin{aligned} \frac{1}{n} \|\mathbf{Z}\boldsymbol{\theta}\|_2^2 &= \boldsymbol{\theta}^\top (\boldsymbol{\Sigma}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta} + \boldsymbol{\theta}^\top \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \boldsymbol{\theta} \\ &\geq \left| \boldsymbol{\theta}^\top (\boldsymbol{\Sigma}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta} \right| - \left| \boldsymbol{\theta}^\top \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \boldsymbol{\theta} \right|. \end{aligned}$$

Note that

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{C}(\mathcal{N}_j) \setminus \{0\}} \frac{\boldsymbol{\theta}^\top (\boldsymbol{\Sigma}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta}}{\|\boldsymbol{\theta}\|_2^2} &\geq \min_{\boldsymbol{\theta} \in \mathcal{M} \setminus \{0\}} \frac{\boldsymbol{\theta}^\top (\boldsymbol{\Sigma}^{\mathbf{X}} \otimes \mathbf{I}_M) \boldsymbol{\theta}}{\|\boldsymbol{\theta}\|_2^2} \\ &\geq \rho_{\min} \left( \boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}} \otimes \mathbf{I}_M \right) \\ &= \rho_{\min}(\boldsymbol{\Sigma}_{\mathcal{N}_j, \mathcal{N}_j}^{\mathbf{X}}) \\ &= \kappa, \end{aligned}$$

and we have  $\kappa(M) > 0$  for all  $M$  by Assumption 4.2. Thus, for any  $\boldsymbol{\theta} \in \mathbb{C}(\mathcal{N}_j)$ , we have

$$\frac{1}{n} \|\mathbf{Z}\boldsymbol{\theta}\|_2^2 \geq \kappa \|\boldsymbol{\theta}\|_2^2 - \left| \boldsymbol{\theta}^\top \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \boldsymbol{\theta} \right|. \quad (89)$$

To prove the RSC condition, it then suffices to give an upper bound for

$$\left| \boldsymbol{\theta}^\top \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \boldsymbol{\theta} \right|,$$

where  $\boldsymbol{\theta} \in \mathbb{C}(\mathcal{N}_j)$ . By Lemma 5 in section D of the appendix of [68] and the definition of  $\mathbb{C}(\mathcal{N}_j)$ , for any  $\boldsymbol{\theta} \in \mathbb{C}(\mathcal{N}_j)$ , we have

$$\begin{aligned} &\left| \boldsymbol{\theta}^\top \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \boldsymbol{\theta} \right| \\ &\leq M^2 \left\| \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \right\|_\infty \|\boldsymbol{\theta}\|_{1,2}^2 \\ &= M^2 \left\| \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \right\|_\infty (\|\boldsymbol{\theta}_{\mathcal{M}}\|_{1,2} + \|\boldsymbol{\theta}_{\mathcal{M}^\perp}\|_{1,2})^2 \\ &\leq 16M^2 \left\| \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \right\|_\infty \|\boldsymbol{\theta}_{\mathcal{M}}\|_{1,2}^2 \\ &\leq 16sM^2 \left\| \left( (\hat{\boldsymbol{\Sigma}}_n^{\mathbf{X}} - \boldsymbol{\Sigma}^{\mathbf{X}}) \otimes \mathbf{I}_M \right) \right\|_\infty \|\boldsymbol{\theta}_{\mathcal{M}}\|_{1,2}^2, \end{aligned}$$

where the penultimate line is by Lemma 6 in section D of [68]. Note that

$$\left\| (\hat{\Sigma}_n^{\mathbf{X}} - \Sigma^{\mathbf{X}}) \otimes \mathbf{I}_M \right\|_{\infty} = \left\| \hat{\Sigma}_n^{\mathbf{X}} - \Sigma^{\mathbf{X}} \right\|_{\infty},$$

then combine (89) and Lemma B.5, we have constant  $c$  such that for any  $\delta > 0$ , with probability at least  $1 - \delta$ , we have

$$\frac{1}{2n} \|\mathbf{Z} \Delta \beta\|_2^2 \geq \left( \frac{\kappa}{2} - 8M^2 s \cdot \left\{ 7\sqrt{3}K_0 \sqrt{\frac{\log(3/\delta) + 2 \log((p-1)M)}{n}} + \frac{8cK_0 \log(2n)(\log(3/\delta) + 2 \log((p-1)M))}{3n} \right\} \right) \|\Delta \beta\|_2^2 \quad (90)$$

for all  $\Delta \beta \in \mathbb{C}(\mathcal{N}_j)$ . Let  $\Gamma(n, p, M, \delta)$  be defined by (88), then when  $\Gamma(n, p, M, \delta) \leq \kappa(M)/(32M^2 s)$ , we have

$$\frac{1}{2n} \|\mathbf{Z} \Delta \beta\|_2^2 \geq \frac{\kappa}{4} \|\Delta \beta\|_2^2 \text{ for all } \Delta \beta \in \mathbb{C}(\mathcal{N}_j)$$

with probability at least  $1 - \delta$ .  $\square$

In the next few results, recall that

$$\mathcal{Q}_{n,\delta}^1 = 1 + 8 \left( \frac{\log(2/\delta)}{n} + \sqrt{\frac{\log(2/\delta)}{n}} \right).$$

**Lemma B.8.** *Let  $\xi_1, \dots, \xi_n$  be i.i.d. mean zero random elements in some Hilbert space, and  $\mathbb{E}[\|\xi_1\|^2] = \sigma_\xi$ . Besides, we assume that*

$$\mathbb{E}[\|\xi_1\|^{2k}] \leq (2\sigma_\xi)^k \cdot k! \quad \text{for all } k = 1, 2, \dots.$$

Then for any given  $\delta \in (0, 1]$ , we have

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|\xi_i\|^2 \leq \sigma_\xi \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta.$$

*Proof.* Note that  $x \mapsto |x|^k$  is a convex function when  $k \geq 1$  and  $x \in \mathbb{R}$ . By Jensen's inequality, we have  $|x/2 + y/2|^k \leq (|x|^k + |y|^k)/2$ , which implies that  $|x + y|^k \leq 2^{k-1}(|x|^k + |y|^k)$ . Then by Lemma B.11, for  $k \geq 2$ , we have

$$\mathbb{E} \left[ \left| \|\xi_1\|^2 - \sigma_\xi \right|^k \right] \leq 2^{k-1} (\mathbb{E}[\|\xi_1\|^{2k}] + \sigma_\xi^k) \leq 2^{k-1} ((2\sigma_\xi)^k k! + \sigma_\xi^k) \leq (4\sigma_\xi)^k k!.$$

Thus, we have

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ \left| \|\xi_i\|^2 - \sigma_\xi \right|^k \right] \leq \frac{k!}{2} (32\sigma_\xi^2)(4\sigma_\xi)^{k-2}$$

for all  $k = 2, 3, \dots$ . Then, by Theorem 2.5 [6], we have

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{i=1}^n \|\xi_i\|^2 - \sigma_\xi \right| > \Delta \right\} \leq 2 \exp \left( -\frac{n\Delta^2}{64\sigma_\xi^2 + 8\sigma_\xi \Delta} \right), \quad \Delta > 0.$$

The result follows by rearranging the terms.  $\square$

**Lemma B.9.** Let  $\delta \in (0, 1]$ . For any  $j \in [p]$ , we have

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|g_{ij}\|^2 \leq \sigma_{\max,0} \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta$$

and

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|e_{ij}\|^2 \leq \sigma_{jr} \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta,$$

where  $e_{ij}$  is defined in (5) and  $\sigma_{jr} = \mathbb{E}[\|e_{ij}\|^2]$ .

*Proof.* The result follows directly from Lemma B.8.  $\square$

**Lemma B.10.** For all  $0 < \delta \leq 1$ , we have

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|\mathbf{r}_i + \mathbf{w}_i\|_2^2 \leq \Xi_4(M) \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta,$$

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|\mathbf{a}_i^{\mathbf{X}^k}\|_2^2 \leq \text{tr}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}}) \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta,$$

and

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n \|a_{im}^{\mathbf{X}^k}\|_2^2 \leq \Xi_2(M) \mathcal{Q}_{n,\delta}^1 \right\} \geq 1 - \delta \text{ for all } m \in [M], k \in [p-1].$$

*Proof.* Note that  $\mathbb{E}[\|\mathbf{r}_i + \mathbf{w}_i\|_2^2] = \Xi_4(M)$ ,  $\mathbb{E}[\|\mathbf{a}_i^{\mathbf{X}^k}\|_2^2] = \text{tr}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}})$  and  $\mathbb{E}[\|a_{im}^{\mathbf{X}^k}\|_2^2] = \Sigma_{mm}^{\mathbf{X}^k} \leq \rho_{\max}(\boldsymbol{\Sigma}_{\mathbf{k}\mathbf{k}}^{\mathbf{X}}) \leq \Xi_2(M)$ , where  $\Xi_2(M), \Xi_4(M)$  is defined in (26), and then the result follows directly from Lemma B.8.  $\square$

**Lemma B.11.** Let  $g$  be a mean zero random Gaussian function in the Hilbert space  $\mathbb{H}$ . Let  $\sigma_0 = \mathbb{E}[\|g\|^2]$ . Then

$$\mathbb{E}[\|g\|^{2k}] \leq (2\sigma_0)^k \cdot k! \quad \text{for all } k = 1, 2, \dots$$

*Proof.* Let  $\{\phi_m\}_{m \geq 1}$  be orthonormal eigenfunctions of covariance function of  $g$ . Let  $a_m = \langle g, \phi_m \rangle$ . We have  $a_m, m \geq 1$ , are independent mean zero Gaussian random variables with variance  $\sigma_m$  and  $\sigma_0 = \sum_{m \geq 1} \sigma_m$ . We further have

$$g = \sum_{m=1}^{\infty} \sigma_m^{1/2} \xi_m \phi_m,$$

where  $\xi_m = \sigma_m^{-1/2} a_m$  are independent standard Gaussian, and

$$\|g\| = \left( \sum_{m \geq 1} \sigma_m \xi_m^2 \right)^{1/2}.$$

Using the Jensen’s inequality for  $t \mapsto t^{2k}$ , we have

$$\begin{aligned} \|g\|^{2k} &= \left( \sum_{m \geq 1} \sigma_m \right)^k \cdot \left( \frac{\sum_{m \geq 1} \sigma_m \xi_m^2}{\sum_{m \geq 1} \sigma_m} \right)^k \\ &\leq \left( \sum_{m \geq 1} \sigma_m \right)^k \cdot \frac{\sum_{m \geq 1} \sigma_m \xi_m^{2k}}{\sum_{m \geq 1} \sigma_m} = \left( \sum_{m \geq 1} \sigma_m \right)^{k-1} \cdot \left( \sum_{m \geq 1} \sigma_m \xi_m^{2k} \right). \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E} [\|g\|^{2k}] &\leq \left( \sum_{m \geq 1} \sigma_m \right)^{k-1} \cdot \left( \sum_{m \geq 1} \sigma_m \mathbb{E} [\xi_m^{2k}] \right) \\ &= \sigma_0^k \mathbb{E} [\xi_1^{2k}] \leq \sigma_0^k \cdot 2^k \cdot k! = (2\sigma_0)^k k!, \end{aligned}$$

which completes the proof. □

**Lemma B.12.** Recall that  $\mathbf{g}_i(\cdot) = (g_{i1}(\cdot), g_{i2}(\cdot), \dots, g_{ip}(\cdot))^\top$  is our  $i$ -th observation defined in Section 2.1. Besides, recall that in Section 2.4, we have  $\phi_m = \phi_{jm}$  and  $\hat{\phi}_m = \hat{\phi}_{jm}$  be the  $m$ -th basis function and its corresponding estimate respectively used to do projection for  $j$ -th node, and  $\hat{\mathbf{a}}_{i,M}^{\mathbf{X}_k} = (\hat{a}_{i1}^{X_k}, \dots, \hat{a}_{iM}^{X_k})^\top$  be the projection score vector of  $\mathbf{g}_i$  by using  $\{\hat{\phi}_m\}_{m=1}^M$ . Under the assumption that

$$\frac{1}{n} \sum_{i=1}^n \|g_{ij}\|^2 \leq \Gamma_g(j) \text{ for all } j \in [p],$$

and

$$\|\hat{\phi}_m - \phi_m\| \leq \Gamma_\phi(m) \text{ for all } m \geq 1,$$

for some  $0 < \delta \leq 1$ , we then have

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{a}_{im}^{X_k} - a_{im}^{X_k} \right)^2 \leq \Gamma_g(k) \Gamma_\phi^2(m),$$

for all  $k \in [p-1]$  and  $m \geq 1$ . Furthermore, we have

$$\frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{a}}_i^{\mathbf{X}_k} - \mathbf{a}_i^{\mathbf{X}_k}\|^2 \leq \Gamma_g(k) \sum_{m=1}^M \Gamma_\phi^2(m).$$

*Proof.* Note that

$$\left( \hat{a}_{im}^{X_k} - a_{im}^{X_k} \right)^2 = \left( \langle g_i^{X_k}, \hat{\phi}_m - \phi_m \rangle \right)^2 \leq \|g_i^{X_k}\|^2 \|\hat{\phi}_m - \phi_m\|^2.$$

Thus we have

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{a}_{im}^{X_k} - a_{im}^{X_k} \right)^2 \leq \left( \frac{1}{n} \sum_{i=1}^n \|g_i^{X_k}\|^2 \right) \|\hat{\phi}_m - \phi_m\|^2.$$

The rest of the proof follows directly from the assumptions we made. □

**Lemma B.13.** For a given node  $j \in [p]$ , recall that  $\phi_m = \phi_{jm}$  and  $\hat{\phi}_m = \hat{\phi}_{jm}$  are the  $m$ -th basis function and its corresponding estimate respectively used to do projection for  $j$ -th node, defined in Section 2.4. Besides, let  $\beta_k(t', t) = \beta_{jk}(t', t)$  and recall that  $b_{k,mm}^* = \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t', t) \phi_m(t) \phi_m(t') dt' dt$  is defined in (7) and  $\tilde{b}_{k,mm} = \int_{\mathcal{T}} \beta_k(t', t) \hat{\phi}_m(t') \hat{\phi}_m(t) dt' dt$  is defined in (44). Under the assumption that

$$\|\hat{\phi}_m - \phi_m\| \leq \Gamma_\phi(m) \text{ for all } m \geq 1,$$

and  $\sum_{m=1}^\infty \Gamma_\phi^2(m) < \infty$ , then we have

$$\begin{aligned} \sum_{k=1}^{p-1} \sum_{m'=1}^\infty (\tilde{b}_{k,mm'} - b_{k,mm'}^*)^2 &\leq 2 \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{HS}^2 \right) \Gamma_\phi^2(m) \\ &\quad + 2 \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^\infty \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) dt' dt \right|^2, \end{aligned}$$

where  $\beta_k(t, t')$  is defined in (5).

*Proof.* Note that

$$\begin{aligned} &|\tilde{b}_{k,mm'} - b_{k,mm'}^*| \\ &= \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') (\hat{\phi}_m(t) \hat{\phi}_{m'}(t') - \phi_m(t) \phi_{m'}(t')) dt' dt \right| \\ &= \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \left\{ (\hat{\phi}_m(t) - \phi_m(t)) \hat{\phi}_{m'}(t') + \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) \right\} dt' dt \right| \\ &\leq \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') (\hat{\phi}_m(t) - \phi_m(t)) \hat{\phi}_{m'}(t') dt' dt \right| \\ &\quad + \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) dt' dt \right|. \end{aligned} \tag{91}$$

Since  $\{\phi_{m'}\}_{m'=1}^\infty$  is an orthonormal function basis, thus when we treat

$$\int_{\mathcal{T}} \beta_k(t, t') (\hat{\phi}_m(t) - \phi_m(t)) dt$$

as a function of  $t'$ , we have

$$\begin{aligned} &\sum_{m'=1}^\infty \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') (\hat{\phi}_m(t) - \phi_m(t)) \hat{\phi}_{m'}(t') dt' dt \right|^2 \\ &= \int_{\mathcal{T}} \left( \int_{\mathcal{T}} \beta_k(t, t') (\hat{\phi}_m(t) - \phi_m(t)) dt \right)^2 dt' \\ &\leq \int_{\mathcal{T}} \left( \int_{\mathcal{T}} \beta_k^2(t, t') dt \right) \left( \int_{\mathcal{T}} (\hat{\phi}_m(t) - \phi_m(t))^2 dt \right) dt' \\ &= \|\hat{\phi}_m - \phi_m\|^2 \|\beta_k(t, t')\|_{HS}^2 \end{aligned}$$



$$\leq \Gamma_\phi^2(m) \|\beta_k(t, t')\|_{HS}^2. \tag{92}$$

Combine (91)-(92), and note that  $b_{k,mm'}^* = \tilde{b}_{k,mm'} = 0$  for all  $m, m' \geq 1$  when  $k \notin \mathcal{N}_j$ , thus we have

$$\begin{aligned} \sum_{k=1}^{p-1} \sum_{m'=1}^{\infty} (\tilde{b}_{k,mm'} - b_{k,mm'}^*)^2 &\leq 2 \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{HS}^2 \right) \Gamma_\phi^2(m) \\ &+ 2 \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) dt' dt \right|^2. \end{aligned}$$

□

**Lemma B.14.** Recall that  $\mathbf{g}_i(\cdot) = (g_{i1}(\cdot), g_{i2}(\cdot), \dots, g_{ip}(\cdot))^\top$  is our  $i$ -th observation defined in Section 2.1 and  $e_{ij}(\cdot)$  is the error term defined in (5). Besides, recall that  $\phi_m = \phi_{jm}$  and  $\hat{\phi}_m = \hat{\phi}_{jm}$  are the  $m$ -th basis function and its corresponding estimate respectively used to do projection for the  $j$ -th node, defined in Section 2.4. Recall from (45)

$$\mathbf{v}_{iM} = \sum_{k=1}^{p-1} (\tilde{\mathbf{B}}_{\mathbf{k},M} - \mathbf{B}_{\mathbf{k},M}^*) \hat{\mathbf{a}}_{i,M}^{X_k} + (\tilde{\mathbf{r}}_{i,M} - \mathbf{r}_{i,M}) + (\tilde{\mathbf{w}}_{i,M} - \mathbf{w}_{i,M}),$$

and suppose that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|g_{ij}\|^2 &\leq \Gamma_g(j) \text{ for all } j \in [p], \\ \frac{1}{n} \sum_{i=1}^n \|e_{ij}\|^2 &\leq \Gamma_e \text{ for all } j \in [p], \end{aligned}$$

and

$$\|\hat{\phi}_m - \phi_m\| \leq \Gamma_\phi(m) \text{ for all } m \geq 1,$$

where  $\sum_{m=1}^\infty \Gamma_\phi^2(m) < \infty$ . Then,

$$\frac{1}{n} \sum_{i=1}^n v_{im}^2 \leq I_m + II_m + III_m, \tag{93}$$

where

$$\begin{aligned} I_m &= 6 \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \left\{ \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{HS}^2 \right) \Gamma_\phi^2(m) \right. \\ &\quad \left. + \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) dt' dt \right|^2 \right\}. \end{aligned}$$

$$II_m = 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \left( \sum_{m'=M+1}^{\infty} \Gamma_{\phi}^2(m') \right),$$

$$III_m = 3\Gamma_e \Gamma_{\phi}^2(m).$$

Recall that  $\beta_{jk}(t, t')$  is defined in (5). Drop the subscript  $j$  and let  $\beta_k(t, t') = \beta_{jk}(t, t')$ . Recall that  $b_{k,mm}^* = \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t', t) \phi_m(t) \phi_{m'}(t') dt' dt$  is defined in (7). Furthermore, we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i\|^2 &\leq 12 \left( \sum_{m=1}^{\infty} \Gamma_{\phi}^2(m) \right) \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{HS}^2 \right) \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \\ &\quad + 3\Phi^2(M) \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \left( \sum_{m'=M+1}^{\infty} \Gamma_{\phi}^2(m') \right) + 3\Gamma_e \sum_{m=1}^M \Gamma_{\phi}^2(m), \end{aligned} \tag{94}$$

where, as previously defined in (37),

$$\Phi(M) = \sqrt{\sum_{k=1}^{p-1} \sum_{m=1}^M \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2}.$$

*Proof.* Note that

$$\begin{aligned} v_{im} &= \sum_{k=1}^{p-1} \sum_{m'=1}^M (\tilde{b}_{k,mm'} - b_{k,mm'}^*) \hat{a}_{im'}^{X_k} \\ &\quad + \sum_{k=1}^{p-1} \sum_{m'=M+1}^{\infty} (\tilde{b}_{k,mm'} \hat{a}_{im'}^{X_k} - b_{k,mm'}^* a_{im'}^{X_k}) + \langle e_{ij}(t), \hat{\phi}_m(t) - \phi_m(t) \rangle \\ &= \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^M (\tilde{b}_{k,mm'} - b_{k,mm'}^*) \hat{a}_{im'}^{X_k} \\ &\quad + \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (\tilde{b}_{k,mm'} \hat{a}_{im'}^{X_k} - b_{k,mm'}^* a_{im'}^{X_k}) + \langle e_{ij}(t), \hat{\phi}_m(t) - \phi_m(t) \rangle \\ &= \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} (\tilde{b}_{k,mm'} - b_{k,mm'}^*) \hat{a}_{im'}^{X_k} + \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} b_{k,mm'}^* (\hat{a}_{im'}^{X_k} - a_{im'}^{X_k}) \\ &\quad + \langle e_{ij}(t), \hat{\phi}_m(t) - \phi_m(t) \rangle. \end{aligned}$$

By Jensen's inequality, we have

$$v_{im}^2 \leq 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} (\tilde{b}_{k,mm'} - b_{k,mm'}^*) \hat{a}_{im'}^{X_k} \right)^2$$

$$\begin{aligned}
 &+ 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} b_{k,mm'}^* (\hat{a}_{im'}^{X_k} - a_{im'}^{X_k}) \right)^2 \\
 &+ 3 \left( (e_{ij}(t), \hat{\phi}_m(t) - \phi_m(t)) \right)^2.
 \end{aligned}$$

By Cauchy-Schwartz inequality, we further have

$$\begin{aligned}
 v_{im}^2 &\leq 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} (\tilde{b}_{k,mm'} - b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} (\hat{a}_{im'}^{X_k})^2 \right) \\
 &+ 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (\hat{a}_{im'}^{X_k} - a_{im'}^{X_k})^2 \right) \\
 &+ 3 \|e_{ij}(t)\|^2 \|\hat{\phi}_m(t) - \phi_m(t)\|^2.
 \end{aligned}$$

Note that  $\sum_{m=1}^{\infty} (\hat{a}_{im}^{X_k})^2 = \|g_i^{X_k}\|^2$ , thus we have

$$\begin{aligned}
 \frac{1}{n} \sum_{i=1}^n v_{im}^2 &\leq 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} (\tilde{b}_{k,mm'} - b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \frac{1}{n} \sum_{i=1}^n \|g_i^{X_k}\|^2 \right) \\
 &+ 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} \frac{1}{n} \sum_{i=1}^n (\hat{a}_{im'}^{X_k} - a_{im'}^{X_k})^2 \right) \\
 &+ 3 \left( \frac{1}{n} \sum_{i=1}^n \|e_{ij}(t)\|^2 \right) \|\hat{\phi}_m(t) - \phi_m(t)\|^2 \\
 &\triangleq \text{I}'_m + \text{II}'_m + \text{III}'_m.
 \end{aligned}$$

By Lemma B.13 and our assumption, we have

$$\begin{aligned}
 \text{I}'_m &\leq 6 \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \left\{ \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{\text{HS}}^2 \right) \Gamma_{\phi}^2(m) \right. \\
 &\quad \left. + \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^{\infty} \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) (\hat{\phi}_{m'}(t') - \phi_{m'}(t')) dt' dt \right|^2 \right\}.
 \end{aligned}$$

By Lemma B.12, we have

$$\frac{1}{n} \sum_{i=1}^n (\hat{a}_{im'}^{X_k} - a_{im'}^{X_k})^2 \leq \Gamma_g(k) \Gamma_{\phi}^2(m').$$

Thus,

$$\text{II}'_m \leq 3 \left( \sum_{k \in \mathcal{N}_j} \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2 \right) \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right) \left( \sum_{m'=M+1}^{\infty} \Gamma_{\phi}^2(m') \right).$$

In addition, we have

$$\text{III}'_m \leq 3\Gamma_\epsilon \Gamma_\phi^2(m).$$

Combining the above results, we complete the proof of (93).

To show (93), note that  $\{\phi_m\}_{m=1}^\infty$  is an orthonormal function basis, and we treat

$$\int_{\mathcal{T}} \beta_k(t, t') \left( \hat{\phi}_{m'}(t') - \phi_{m'}(t') \right) dt'$$

as a function of  $t$ , then we have

$$\begin{aligned} & \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^\infty \sum_{m=1}^\infty \left| \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t, t') \phi_m(t) \left( \hat{\phi}_{m'}(t') - \phi_{m'}(t') \right) dt' dt \right|^2 \\ &= \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^\infty \int_{\mathcal{T}} \left\{ \int_{\mathcal{T}} \beta_k(t, t') \left( \hat{\phi}_{m'}(t') - \phi_{m'}(t') \right) dt' \right\}^2 dt \\ &\leq \sum_{k \in \mathcal{N}_j} \sum_{m'=1}^\infty \int_{\mathcal{T}} \left\{ \int_{\mathcal{T}} \beta_k^2(t, t') dt \right\} \left\{ \int_{\mathcal{T}} \left( \hat{\phi}_{m'}(t') - \phi_{m'}(t') \right)^2 dt' \right\} dt \\ &= \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{\text{HS}}^2 \right) \left( \sum_{m'=1}^\infty \Gamma_\phi^2(m') \right). \end{aligned}$$

Thus, we have

$$\sum_{m=1}^M I_m \leq \sum_{m=1}^\infty I_m \leq 12 \left( \sum_{m=1}^\infty \Gamma_\phi^2(m) \right) \left( \sum_{k \in \mathcal{N}_j} \|\beta_k(t, t')\|_{\text{HS}}^2 \right) \left( \sum_{k \in \mathcal{N}_j} \Gamma_g(k) \right).$$

(94) then follows the combination of the above inequality and (93). □

**Lemma B.15.** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ . We have

$$\sqrt{\rho_{\min}(\mathbf{B}^\top \mathbf{B})} \|\mathbf{A}\|_F \leq \|\mathbf{BA}\|_F \leq \|\mathbf{B}\|_2 \|\mathbf{A}\|_F$$

and

$$\sqrt{\rho_{\min}(\mathbf{BB}^\top)} \|\mathbf{A}\|_F \leq \|\mathbf{AB}\|_F \leq \|\mathbf{B}\|_2 \|\mathbf{A}\|_F.$$

*Proof.* Let  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$ . Since

$$\|\mathbf{BA}\|_F^2 = \text{tr}(\mathbf{A}^\top \mathbf{B}^\top \mathbf{BA}) = \sum_{i=1}^n \mathbf{a}_i^\top \mathbf{B}^\top \mathbf{B} \mathbf{a}_i,$$

we have

$$\rho_{\min}(\mathbf{B}^\top \mathbf{B}) \|\mathbf{A}\|_F^2 = \rho_{\min}(\mathbf{B}^\top \mathbf{B}) \sum_{i=1}^n \|\mathbf{a}_i\|_2^2$$

$$\begin{aligned}
&\leq \sum_{i=1}^n \mathbf{a}_i^\top \mathbf{B}^\top \mathbf{B} \mathbf{a}_i \\
&\leq \|\mathbf{B}\|_2^2 \sum_{i=1}^n \|\mathbf{a}_i\|_2^2 \\
&= \|\mathbf{B}\|_2^2 \|\mathbf{A}\|_F^2.
\end{aligned}$$

The final result follows from taking the square root on both sides.  $\square$

**Lemma B.16.** *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric and positive semi-definite,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ , and  $\mathbf{C} \in \mathbb{R}^{n \times n}$  be a diagonal matrix with positive diagonal elements. Then we have*

$$\|(\mathbf{A} + \mathbf{C})\mathbf{B}\|_F \geq \|\mathbf{C}\mathbf{B}\|_F.$$

*Proof.* Let  $j$ -th column of  $\mathbf{B}$  be  $\mathbf{b}_j$  for  $j = 1, \dots, m$ . Then we have

$$\begin{aligned}
\|(\mathbf{A} + \mathbf{C})\mathbf{B}\|_F^2 &= \text{tr}(\mathbf{B}^\top (\mathbf{A} + \mathbf{C})^\top (\mathbf{A} + \mathbf{C}) \mathbf{B}) \\
&= \sum_{j=1}^m \mathbf{b}_j^\top (\mathbf{A} + \mathbf{C})^\top (\mathbf{A} + \mathbf{C}) \mathbf{b}_j \\
&= \sum_{j=1}^m \mathbf{b}_j^\top \mathbf{A}^\top \mathbf{A} \mathbf{b}_j + \sum_{j=1}^m \mathbf{b}_j^\top \mathbf{C}^\top \mathbf{A} \mathbf{b}_j \\
&\quad + \sum_{j=1}^m \mathbf{b}_j^\top \mathbf{A}^\top \mathbf{C} \mathbf{b}_j + \sum_{j=1}^m \mathbf{b}_j^\top \mathbf{C}^\top \mathbf{C} \mathbf{b}_j.
\end{aligned}$$

Since  $\mathbf{A}$  is symmetric and positive semi-definite and  $\mathbf{C}$  is a diagonal matrix with positive diagonal elements, we have  $\mathbf{C}^\top \mathbf{A} = \mathbf{A}^\top \mathbf{C} = \mathbf{C} \mathbf{A}$  to be symmetric and positive semi-definite. And note that  $\mathbf{A}^\top \mathbf{A}$  is symmetric and positive semi-definite, thus we have

$$\|(\mathbf{A} + \mathbf{C})\mathbf{B}\|_F^2 \geq \sum_{j=1}^m \mathbf{b}_j^\top \mathbf{C}^\top \mathbf{C} \mathbf{b}_j = \text{tr}(\mathbf{B}^\top \mathbf{C}^\top \mathbf{C} \mathbf{B}) = \|\mathbf{C}\mathbf{B}\|_F^2,$$

which implies the final result.  $\square$

### Appendix C: Wall-clock runtime comparison

Table 4 documents the wall-clock runtime required for each method to generate a Receiver Operating Characteristic (ROC) curve. We execute each method on 100 distinct values of  $\lambda_n$ . Table 5 logs the wall-clock runtime of the FPCA- $g_X$  method to generate an estimated graph via the Selective Cross-Validation (SCV) algorithm as described in Algorithm 2, over a two-dimensional grid of  $(\lambda_n, t_\epsilon)$ . This grid encompasses 100 unique values of  $\lambda_n$  and 7 different values of  $t_\epsilon$ . All tasks are executed in parallel using R, utilizing 28 CPU cores on Chicago Booth's Mercury Computing Cluster.

TABLE 4

The average wall-clock running time for each method to obtain a ROC under each model. All of results are recorded in seconds. The fastest algorithm in each setting is marked bold. The standard deviation is given in parenthesis.

Model	$p$	FPCA- $g_X$	FPCA- $g_Y$	FGLasso	PSKL	FPCA-PSKL
A	50	24.8 (3.0)	47.9 (6.3)	<b>11.4</b> (2.5)	126.6 (9.9)	N/A
	100	137.7 (16.7)	186.4 (20.7)	<b>35.1</b> (21.4)	241.3 (21.2)	
	150	305.9 (58.2)	377.9 (63.8)	<b>69.9</b> (10.3)	418.4 (24.0)	
B	50	46.6 (2.6)	75.1 (3.9)	<b>33.5</b> (15.5)	108.1 (7.8)	N/A
	100	263.3 (29.2)	345.0 (33.4)	<b>111.0</b> (72.1)	226.8 (13.4)	
	150	807.7 (125.3)	973.7 (136.0)	<b>212.5</b> (188.2)	437.3 (27.9)	
C	50	<b>14.6</b> (4.1)	29.2 (7.1)	91.3 (27.7)	343.3 (43.5)	N/A
	100	174.1 (125.2)	218.9 (113.1)	<b>87.8</b> (78.8)	695.6 (73.3)	
	150	692.3 (684.0)	754.9 (517.1)	<b>139.8</b> (132.7)	1037.9 (200.2)	
D	50	<b>11.2</b> (16.6)	26.7 (22.3)	72.3 (37.0)	259.6 (74.8)	14.8 (23.0)
	100	<b>92.0</b> (99.1)	129.9 (121.4)	390.2 (333.8)	603.3 (151.0)	112.3 (121.8)
	150	<b>92.3</b> (111.5)	126 (128.2)	1228.2 (325.7)	1108.4 (103.6)	100.0 (122.3)

TABLE 5

The average wall-clock running time of FPCA- $g_X$  to obtain an estimated graph under each model, where  $(\lambda_n, \epsilon_n)$  is chosen by SCV algorithm stated in Algorithm 2. All of results are recorded in seconds. The standard deviation is given in parenthesis.

Model	A	B	C	D
$p = 50$	22.3 (2.9)	45.0 (6.9)	12.2 (1.2)	17.0 (16.2)
$p = 100$	141.1 (24.1)	218.9 (42.8)	111.7 (12.9)	108.6 (106.8)
$p = 150$	326.0 (70.5)	707.1 (174.0)	463.8 (322.0)	220.6 (60.0)

The runtime analysis indicates that the FPCA- $g_X$  method is marginally faster than the other two FPCA methods and significantly outpaces PSKL in most instances. Even though FPCA- $g_X$  can occasionally be slower than FGLasso, it delivers more accurate results while maintaining comparable runtime. Notably, FPCA methods outperform FGLasso or PSKL in terms of speed. Despite the SCV process operating over a two-dimensional grid with 100 different values of  $\lambda_n$  and 7 different values of  $t_\epsilon$ , its runtime is akin to that of the ROC process, which operates on a one-dimensional grid of  $\lambda_n$ . This can be attributed to the fact that the most time-intensive step of the SCV process, the ADMM algorithm, is executed only once for each value of  $\lambda_n$ .

**Appendix D: Labels of ROIs in the AAL atlas**

Table 6: Labels and names of each ROI in the AAL atlas

Label	Name	Label	Name
2001	Precentral_L	4112	ParaHippocampal_R
2002	Precentral_R	4201	Amygdala_L
2101	Frontal_Sup_L	4202	Amygdala_R
2102	Frontal_Sup_R	5001	Calcarine_L
2111	Frontal_Sup_Orb_L	5002	Calcarine_R
2112	Frontal_Sup_Orb_R	5011	Cuneus_L
2201	Frontal_Mid_L	5012	Cuneus_R
2202	Frontal_Mid_R	5021	Lingual_L
2211	Frontal_Mid_Orb_L	5022	Lingual_R
2212	Frontal_Mid_Orb_R	5101	Occipital_Sup_L
2301	Frontal_Inf_Oper_L	5102	Occipital_Sup_R
2302	Frontal_Inf_Oper_R	5201	Occipital_Mid_L
2311	Frontal_Inf_Tri_L	5202	Occipital_Mid_R
2312	Frontal_Inf_Tri_R	5301	Occipital_Inf_L
2321	Frontal_Inf_Orb_L	5302	Occipital_Inf_R
2322	Frontal_Inf_Orb_R	5401	Fusiform_L
2331	Rolandic_Oper_L	5402	Fusiform_R
2332	Rolandic_Oper_R	6001	Postcentral_L
2401	Supp_Motor_Area_L	6002	Postcentral_R
2402	Supp_Motor_Area_R	6101	Parietal_Sup_L
2501	Olfactory_L	6102	Parietal_Sup_R
2502	Olfactory_R	6201	Parietal_Inf_L
2601	Frontal_Sup_Medial_L	6202	Parietal_Inf_R
2602	Frontal_Sup_Medial_R	6211	SupraMarginal_L
2611	Frontal_Med_Orb_L	6212	SupraMarginal_R
2612	Frontal_Med_Orb_R	6221	Angular_L
2701	Rectus_L	6222	Angular_R
2702	Rectus_R	6301	Precuneus_L
3001	Insula_L	6302	Precuneus_R
3002	Insula_R	6401	Paracentral_Lobule_L
4001	Cingulum_Ant_L	6402	Paracentral_Lobule_R
4002	Cingulum_Ant_R	7001	Caudate_L
4011	Cingulum_Mid_L	7002	Caudate_R
4012	Cingulum_Mid_R	7011	Putamen_L
4021	Cingulum_Post_L	7012	Putamen_R
4022	Cingulum_Post_R	7021	Pallidum_L
4101	Hippocampus_L	7022	Pallidum_R
4102	Hippocampus_R	7101	Thalamus_L
4111	ParaHippocampal_L	7102	Thalamus_R
8101	Heschl_L	8111	Temporal_Sup_L

8102	Heschl_R	8112	Temporal_Sup_R
8121	Temporal_Pole_Sup_L		
8122	Temporal_Pole_Sup_R		
8201	Temporal_Mid_L		
8202	Temporal_Mid_R		
8211	Temporal_Pole_Mid_L		
8212	Temporal_Pole_Mid_R		
8301	Temporal_Inf_L		
8302	Temporal_Inf_R		
9001	Cerebelum_Crus1_L		
9002	Cerebelum_Crus1_R		
9011	Cerebelum_Crus2_L		
9012	Cerebelum_Crus2_R		
9021	Cerebelum_3_L		
9022	Cerebelum_3_R		
9031	Cerebelum_4_5_L		
9032	Cerebelum_4_5_R		
9041	Cerebelum_6_L		
9042	Cerebelum_6_R		
9051	Cerebelum_7b_L		
9052	Cerebelum_7b_R		
9061	Cerebelum_8_L		
9062	Cerebelum_8_R		
9071	Cerebelum_9_L		
9072	Cerebelum_9_R		
9081	Cerebelum_10_L		
9082	Cerebelum_10_R		
9100	Vermis_1_2		
9110	Vermis_3		
9120	Vermis_4_5		
9130	Vermis_6		
9140	Vermis_7		
9150	Vermis_8		
9160	Vermis_9		
9170	Vermis_10		



## Appendix E: Table of notations

Table 7: Summary of notations used in the paper.

Notation	Meaning	Page
$G = (V, E)$	undirected graph, $V$ is set of vertices, $E$ is set of edges	2
$\mathbf{X}$	$p$ -dimensional random variables	2
$M$	number of basis functions we used to do dimension reduction	3
$\mathbf{g}(\cdot)$	$p$ -dimensional multivariate Gaussian process	5
$\mathcal{T}$	domain of multivariate Gaussian process	5
$C_{jl}(t, t')$	conditional cross-covariance function	5
$\{\mathbf{g}_i(\cdot)\}_{i=1}^n$	$\mathbf{g}_i(\cdot) = (g_{i1}(\cdot), \dots, g_{ip}(\cdot))^\top$ , random copies of $\mathbf{g}(\cdot)$	5
$\beta_{jk}(t, t')$	coefficient on $g_{ij}$ from $g_{ik}$	5
$\mathcal{N}_j, \hat{\mathcal{N}}_j$	(estimated) neighborhood set of node $j$	5
$e_{ij}(\cdot)$	error of $g_{ij}(\cdot)$	5
$\phi_j$	$= \{\phi_{jm}(\cdot)\}_{m=1}^\infty$ , orthonormal functional basis on $\mathbb{H}$	6
$\mathbf{a}_{i,k,M}$	$= (a_{ik1}, \dots, a_{ikM})^\top$ , vector of projection scores	6
$g_i^Y(\cdot), g_i^{X_k}(\cdot)$	random functions of the target node and the other random functions	6
$\mathbf{a}_{i,M}^Y, \mathbf{a}_{i,M}^{X_k}, \mathbf{a}_{i,M}^X$	vectors of scores projected on known bases $\phi_j$	6
$\mathbf{B}_{k,M}^*$	regression matrix parameter	7
$\mathbf{w}_{i,M}, \mathbf{r}_{i,M}$	noise vector and bias term for $M$ -truncation	7
$\hat{\mathbf{B}}_{k,M}$	Estimator of $\mathbf{B}_{k,M}^*$	7
$\lambda_n$	penalty parameter for group Lasso regression	7
$\epsilon_n$	threshold parameter of neighborhood recognition	7
$\hat{E}$	estimated edge set	7
$\hat{\mathbf{a}}_{i,M}^Y, \hat{\mathbf{a}}_{i,M}^{X_k}, \hat{\mathbf{a}}_{i,M}^X$	scores projected on estimated bases $\hat{\phi}_j$	7-8
$K_{jj}(t', t)$	functional covariance of $g_{ij}(\cdot)$	8
$\mathcal{K}_j(f)(t)$	Hilbert-Schmidt covariance operator	8
$\{\sigma_{jm}\}_{m \in \mathbb{N}}$	eigenvalues of $\mathcal{K}_j$	8
$g_{ij}^M(t)$	$\mathcal{L}^2$ projection of $g_{ij}(t)$ onto the basis spanned by the first $M$ FPCA functions	8
$\hat{K}_{jj}(t', t)$	empirical functional covariance of $g_{ij}(\cdot)$	8
$\{\hat{\sigma}_{jm}, \hat{\phi}_{jm}(t)\}_{m=1}^M$	eigenpairs of $\hat{K}_{jj}(t', t)$	8
$\hat{a}_{ijm}$	$= \int_{\mathcal{T}} g_{ij}(t) \hat{\phi}_{jm}(t) dt$ , estimated FPCA scores	8
$\hat{\mathbf{B}}_{\lambda_n}$	$= (\hat{\mathbf{B}}_1, \dots, \hat{\mathbf{B}}_{p-1})$ , group Lasso estimates under a fixed $\lambda_n$	9

$\tilde{\mathbf{B}}_k$	estimates of $\mathbf{B}_k$ in selective cross-validation process	10
$\hat{\boldsymbol{\epsilon}}_i$	$= \mathbf{a}_{i,M}^Y - \sum_{k=1}^{p-1} \tilde{\mathbf{B}}_k \mathbf{a}_{i,M}^{\mathbf{X}_k}$ , residuals of $\tilde{\mathbf{B}}_k$	10
$\mathbf{A}^Y, \mathbf{A}^{\mathbf{X}_k}, \mathbf{A}^{\mathbf{X}}$	matrices of FPCA scores	10
$\rho, \rho^h$	penalty parameter for ADMM subproblem	10
$b_{k,mm}^*$	$= \int_{\mathcal{T} \times \mathcal{T}} \beta_k(t', t) \phi_m(t) \phi_{m'}(t') dt' dt$	12
$\beta_{k,M}(t', t)$	$= \sum_{m,m'=1}^M b_{k,mm'}^* \phi_m(t) \phi_{m'}(t')$	13
$\beta_{k,>M}(t', t)$	$= \sum_{m>M \text{ or } m'>M}^{\infty} b_{k,mm'}^* \phi_m(t) \phi_{m'}(t')$	13
$\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{r}}$	$= \text{Cov}(\mathbf{a}_{i,M}^{\mathbf{X}_k}, \mathbf{r}_{i,M})$ , and $\boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{X}_k} = (\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{r}})^\top$	13
$\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{X}_l}$	$= \text{Cov}(\mathbf{a}_{i,M}^{\mathbf{X}_k}, \mathbf{a}_{i,M}^{\mathbf{X}_l})$	13
$\Xi_1(M)$	$= \max_{k \in [p-1]} \{ \rho_{\max}(\boldsymbol{\Sigma}^{\mathbf{w}} + \boldsymbol{\Sigma}^{\mathbf{r}} - \boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{X}_k} (\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{X}_k})^{-1} \boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{r}}) \}$	13
$\Xi_2(M)$	$= \max_{k \in [p-1]} \rho_{\max}(\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{X}_k})$	13
$\Xi_3(M)$	$= \max_{k \in [p-1]} \text{tr}(\boldsymbol{\Sigma}^{\mathbf{X}_k, \mathbf{X}_k})$	13
$\Xi_4(M)$	$= \text{tr} \{ \boldsymbol{\Sigma}^{\mathbf{r}} + \boldsymbol{\Sigma}^{\mathbf{w}} + \boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{w}} + \boldsymbol{\Sigma}^{\mathbf{w}, \mathbf{r}} \}$	13
$\omega(M)$	$= \max_{k \in [p-1]} \ \boldsymbol{\Sigma}^{\mathbf{r}, \mathbf{X}_k}\ _{\text{F}}$	13
$K_0$	$= \max_{k \in [p-1], m \in M} \mathbb{E}[(a_{i,m}^{\mathbf{X}_k})^2]$ , used to derive an upper bound for the estimation error of the covariance matrix of $\mathbf{a}_{i,M}^{\mathbf{X}}$	13
$\kappa(M)$	$= \rho_{\min}((\boldsymbol{\Sigma}^{\mathbf{X}})_{\mathcal{N}_j, \mathcal{N}_j})$	14
$\tau(M)$	$= \min_{k \in \mathcal{N}_j} \ \mathbf{B}_k^*\ _{\text{F}} = \min_{k \in \mathcal{N}_j} \ \beta_{k,M}(t', t)\ _{\text{HS}}$ , relevant signal strength	14
$\chi(n, p, M, \delta)$	$= \frac{6\sqrt{s}}{\sqrt{\kappa(M)}} \tilde{\lambda}(n, p, M, \delta)$	14
$\tilde{\lambda}(n, p, M, \delta)$	exact form can be found in (54)	14
$\Phi(M)$	$= \sqrt{\sum_{k=1}^{p-1} \sum_{m=1}^M \sum_{m'=M+1}^{\infty} (b_{k,mm'}^*)^2}$	16
$\check{\chi}(n, p, M, \delta)$	$= \frac{6\sqrt{s}}{\sqrt{\kappa(M)}} \check{\lambda}(n, p, M, \delta)$	16
$\check{\lambda}(n, p, M, \delta)$	exact form given in (69)	16
$\Lambda(M, \phi)$	$\frac{\omega(M)}{\sqrt{\kappa(M)}\tau(M)}$	17
$(\mathbf{a}_{i1}, \dots, \mathbf{a}_{ip})^\top$	generated functional scores from a mean zero Gaussian distribution	18
$\mathbf{f}(\cdot)$	vector of Fourier basis functions	18
$M^*$	number of basis functions contained in $\mathbf{f}(\cdot)$	18
$\boldsymbol{\Sigma}$	covariance matrix of $(\mathbf{a}_{i1}, \dots, \mathbf{a}_{ip})^\top$	18
$\boldsymbol{\Theta}$	$= \boldsymbol{\Sigma}^{-1}$ , precision matrix	18
$T$	number of observation time points	19
$\epsilon_{ijk}$	observation error of $g_{ij}(t_k)$	19
$\sigma$	variance of $\epsilon_{ijk}$	19

## Acknowledgments

The authors would like to thank the anonymous referees, an Associate Editor and the Editor for their constructive comments that improved the quality of this paper. We would also like to thank Zhaohan Wu from Florida State University for his suggestions on fMRI data analysis. This work was completed in part with resources provided by the University of Chicago Booth Mercury Computing Cluster.

## Funding

The research of MK is supported in part by NSF Grant ECCS-2216912.

## Supplementary Material

### FGM Neighborhood Code

(doi: [10.1214/24-EJS2219SUPP](https://doi.org/10.1214/24-EJS2219SUPP); .zip). This file contains the code and data used for the experiments in Section 5 and Section 6

## References

- [1] ADLER, R. J. and TAYLOR, J. E. (2007). *Random Fields and Geometry*. Springer, New York. [MR2319516](#)
- [2] ALLEN, G., MÜLLER, R.-A. and COURCHESNE, E. (2004). Cerebellar function in autism: functional magnetic resonance image activation during a simple motor task. *Biological psychiatry* **56** 269–278.
- [3] BELLEC, P., CHU, C., CHOUINARD-DECORTE, F., BENHAJALI, Y., MARGULIES, D. S. and CRADDOCK, R. C. (2017). The Neuro Bureau ADHD-200 Preprocessed Repository. *Neuroimage* **144** 275–286.
- [4] BELMONTE, M. K., ALLEN, G., BECKEL-MITCHENER, A., BOULANGER, L. M., CARPER, R. A. and WEBB, S. J. (2004). Autism and abnormal development of brain connectivity. *Journal of Neuroscience* **24** 9228–9231.
- [5] BESAG, J. (1975). Statistical Analysis of Non-Lattice Data. *Journal of the Royal Statistical Society. Series D* **24** 179–195.
- [6] BOSQ, D. (2000). *Linear Processes in Function Spaces*. Springer-Verlag, New York. [MR1783138](#)
- [7] BOYD, S. P., PARIKH, N., CHU, E., PELEATO, B. and ECKSTEIN, J. (2011). Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends in Machine Learning* **3** 1–122.
- [8] BU, X., CAO, M., HUANG, X. and HE, Y. (2021). The structural connectome in ADHD. *Psychoradiology* **1** 257–271.
- [9] BUBECK, S. (2015). Convex Optimization: Algorithms and Complexity. *Foundations and Trends in Machine Learning* **8** 231–357.

- [10] CAI, T., LIU, W. and LUO, X. (2011). A Constrained L1 Minimization Approach to Sparse Precision Matrix Estimation. *Journal of the American Statistical Association* **106** 594-607. [MR2847973](#)
- [11] CAO, Q., SHU, N., AN, L., WANG, P., SUN, L., XIA, M.-R., WANG, J.-H., GONG, G.-L., ZANG, Y.-F., WANG, Y.-F. et al. (2013). Probabilistic diffusion tractography and graph theory analysis reveal abnormal white matter structural connectivity networks in drug-naive boys with attention deficit/hyperactivity disorder. *Journal of Neuroscience* **33** 10676–10687.
- [12] CARIA, A., CIRINGIONE, L. and DE FALCO, S. (2020). Morphofunctional alterations of the hypothalamus and social behavior in autism spectrum disorders. *Brain Sciences* **10** 435.
- [13] CHIOU, J. M. and MÜLLER, H. G. (2016). A pairwise interaction model for multivariate functional and longitudinal data. *Biometrika* **103** 377–396. [MR3509893](#)
- [14] CHRISTENSEN, D. L., BRAUN, K. V. N., BAIO, J., BILDER, D., CHARLES, J., CONSTANTINO, J. N., DANIELS, J., DURKIN, M. S., FITZGERALD, R. T., KURZIUS-SPENCER, M. et al. (2018). Prevalence and characteristics of autism spectrum disorder among children aged 8 years—autism and developmental disabilities monitoring network, 11 sites, United States, 2012. *MMWR Surveillance Summaries* **65** 1.
- [15] CRADDOCK, C., BENHAJALI, Y., CHU, C., CHOUINARD, F., EVANS, A., JAKAB, A., KHUNDRAPAM, B. S., LEWIS, J. D., LI, Q., MILHAM, M. et al. (2013). The neuro bureau preprocessing initiative: open sharing of preprocessed neuroimaging data and derivatives. *Frontiers in Neuroinformatics* **7** 27. [MR3855624](#)
- [16] DI MARTINO, A., YAN, C.-G., LI, Q., DENIO, E., CASTELLANOS, F. X., ALAERTS, K., ANDERSON, J. S., ASSAF, M., BOOKHEIMER, S. Y., DAPRETTO, M. et al. (2014). The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular psychiatry* **19** 659–667.
- [17] GABAY, D. and MERCIER, B. (1976). A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & mathematics with applications* **2** 17–40.
- [18] GLICKSTEIN, M., STRATA, P. and VOOGD, J. (2009). Cerebellum: history. *Neuroscience* **162** 549–559.
- [19] GRILL-SPECTOR, K. and MALACH, R. (2004). The human visual cortex. *Annual Review of Neuroscience* **27** 649–677.
- [20] HE, B., YANG, H. and WANG, S. (2000). Alternating direction method with self-adaptive penalty parameters for monotone variational inequalities. *Journal of Optimization Theory and applications* **106** 337–356. [MR1788928](#)
- [21] HONG, S.-J., VOS DE WAEL, R., BETHLEHEM, R. A., LARIVIERE, S., PAQUOLA, C., VALK, S. L., MILHAM, M. P., DI MARTINO, A., MARGULIES, D. S., SMALLWOOD, J. et al. (2019). Atypical functional connectome hierarchy in autism. *Nature communications* **10** 1022.
- [22] HSING, T. and EUBANK, R. (2015). *Theoretical foundations of functional data analysis, with an introduction to linear operators*. John Wiley & Sons,

- Ltd., Chichester. [MR3379106](#)
- [23] IVANESCU, A. E., STAICU, A. M., SCHEIPL, F. and GREVEN, S. (2015). Penalized function-on-function regression. *Computational Statistics* **30** 539–568. [MR3357075](#)
- [24] JANSON, S. (1997). *Gaussian Hilbert Spaces*. Cambridge University Press. [MR1474726](#)
- [25] KALLENBERG, O. (1997). *Foundations of modern probability. Probability and its Applications*. Springer-Verlag, New York. [MR1464694](#)
- [26] KIM, B., LIU, S. and KOLAR, M. (2021). Two-sample inference for high-dimensional Markov networks. *Journal of the Royal Statistical Society. Series B. Statistical Methodology* **83** 939–962. [MR4349123](#)
- [27] KLEBANOV, I., SPRUNGK, B. and SULLIVAN, T. J. (2021). The linear conditional expectation in Hilbert space. *Bernoulli* **27** 2267–2299. [MR4303883](#)
- [28] KOLAR, M., LIU, H. and XING, E. P. (2013). Markov Network Estimation From Multi-attribute Data. In *International Conference on Machine Learning, ICML*.
- [29] KOLAR, M., LIU, H. and XING, E. P. (2014). Graph estimation from multi-attribute data. *Journal of Machine Learning Research* **15** 1713–1750. [MR3225245](#)
- [30] KONRAD, K., NEUFANG, S., HANISCH, C., FINK, G. R. and HERPERTZ-DAHLMANN, B. (2006). Dysfunctional attentional networks in children with attention deficit/hyperactivity disorder: evidence from an event-related functional magnetic resonance imaging study. *Biological psychiatry* **59** 643–651.
- [31] KUCHIBHOTLA, A. K. and CHAKRABORTTY, A. (2022). Moving beyond sub-Gaussianity in high-dimensional statistics: applications in covariance estimation and linear regression. *Information and Inference* **11** 1389–1456. [MR4526326](#)
- [32] LAURITZEN, S. L. (1996). *Graphical Models. Oxford Statistical Science Series*. Clarendon Press. [MR1419991](#)
- [33] LEE, K.-Y., LI, B. and ZHAO, H. (2016). Variable selection via additive conditional independence. *Journal of the Royal Statistical Society: Series B* **78** 1037–1055. [MR3557188](#)
- [34] LEE, K.-Y., LI, B. and ZHAO, H. (2016). On an additive partial correlation operator and nonparametric estimation of graphical models. *Biometrika* **103** 513–530. [MR3551781](#)
- [35] LEE, K.-Y., LI, L., LI, B. and ZHAO, H. (2022). Nonparametric Functional Graphical Modeling Through Functional Additive Regression Operator. *Journal of the American Statistical Association* **0** 1–15. [MR4646601](#)
- [36] LI, B., CHUN, H. and ZHAO, H. (2014). On an additive semigraphoid model for statistical networks with application to pathway analysis. *Journal of the American Statistical Association* **109** 1188–1204. [MR3265690](#)
- [37] LI, B. and SOLEA, E. (2018). A nonparametric graphical model for functional data with application to brain networks based on fMRI. *Journal of the American Statistical Association* **113** 1637–1655. [MR3902235](#)
- [38] LUO, R. and QI, X. (2017). Function-on-function linear regression by signal

- compression. *Journal of the American Statistical Association* **112** 690–705. [MR3671763](#)
- [39] LUO, R., QI, X. and WANG, Y. (2016). Functional wavelet regression for linear function-on-function models. *Electronic Journal of Statistics* **10** 3179–3216. [MR3571966](#)
- [40] LYNCH, B. and CHEN, K. (2018). A test of weak separability for multi-way functional data, with application to brain connectivity studies. *Biometrika* **105** 815–831. [MR3877867](#)
- [41] MAXIMO, J. O., CADENA, E. J. and KANA, R. K. (2014). The implications of brain connectivity in the neuropsychology of autism. *Neuropsychology review* **24** 16–31.
- [42] MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High Dimensional Graphs And Variable Selection With The Lasso. *Annals of Statistics* **34** 1436–1462. [MR2278363](#)
- [43] MEINSHAUSEN, N. and YU, B. (2009). Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics* **37** 246–270. [MR2488351](#)
- [44] MESULAM, M. (2012). The evolving landscape of human cortical connectivity: facts and inferences. *Neuroimage* **62** 2182–2189.
- [45] MILHAM, M. P., FAIR, D., MENNES, M., MOSTOFISKY, S. H. et al. (2012). The ADHD-200 consortium: a model to advance the translational potential of neuroimaging in clinical neuroscience. *Frontiers in systems neuroscience* **6** 62.
- [46] NEBEL, M. B., ELOYAN, A., BARBER, A. D. and MOSTOFISKY, S. H. (2014). Precentral gyrus functional connectivity signatures of autism. *Frontiers in systems neuroscience* **8** 80.
- [47] NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A unified framework for high-dimensional analysis of  $M$ -estimators with decomposable regularizers. *Statistical Science* **27** 538–557. [MR3025133](#)
- [48] NOORDERMEER, S. D., LUMAN, M., GREVEN, C. U., VEROUDE, K., FARAONE, S. V., HARTMAN, C. A., HOEKSTRA, P. J., FRANKE, B., BUITELAAR, J. K., HESLENFELD, D. J. et al. (2017). Structural brain abnormalities of attention-deficit/hyperactivity disorder with oppositional defiant disorder. *Biological Psychiatry* **82** 642–650.
- [49] PATRIQUIN, M. A., DERAMUS, T., LIBERO, L. E., LAIRD, A. and KANA, R. K. (2016). Neuroanatomical and neurofunctional markers of social cognition in autism spectrum disorder. *Human brain mapping* **37** 3957–3978.
- [50] QI, X. and LUO, R. (2018). Function-on-function regression with thousands of predictive curves. *Journal of Multivariate Analysis* **163** 51–66. [MR3732340](#)
- [51] QI, X. and LUO, R. (2019). Nonlinear function-on-function additive model with multiple predictor curves. *Statistica Sinica* **29** 719–739. [MR3931385](#)
- [52] QIAO, X., GUO, S. and JAMES, G. M. (2019). Functional Graphical Models. *Journal of the American Statistical Association* **114** 211–222.

- [MR3941249](#)
- [53] QIAO, X., QIAN, C., JAMES, G. M. and GUO, S. (2020). Doubly functional graphical models in high dimensions. *Biometrika* **107** 415–431. [MR4108937](#)
- [54] RAMSAY, J. O. and SILVERMAN, B. W. (2005). *Functional data analysis*, Second ed. Springer, New York. [MR2168993](#)
- [55] ROTHMAN, A. J., BICKEL, P. J., LEVINA, E., ZHU, J. et al. (2008). Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics* **2** 494–515. [MR2417391](#)
- [56] SCHEIPL, F., STAIKU, A. M. and GREVEN, S. (2015). Functional additive mixed models. *Journal of Computational and Graphical Statistics* **24** 477–501. [MR3357391](#)
- [57] SHE, Y. (2012). An iterative algorithm for fitting nonconvex penalized generalized linear models with grouped predictors. *Computational Statistics & Data Analysis* **56** 2976–2990. [MR2929353](#)
- [58] SOLEA, E. and DETTE, H. (2022). Nonparametric and high-dimensional functional graphical models. *Electronic Journal of Statistics* **16** 6175–6231. [MR4515718](#)
- [59] SOLEA, E. and LI, B. (2022). Copula Gaussian graphical models for functional data. *Journal of the American Statistical Association* **117** 781–793. [MR4436312](#)
- [60] TSAI, K., ZHAO, B., KOYEJO, S. and KOLAR, M. (2023). Latent Multimodal Functional Graphical Model Estimation. *Journal of the American Statistical Association* **0** 1–25. [MR4515044](#)
- [61] TSAY, R. S. and POURAHMADI, M. (2017). Modelling structured correlation matrices. *Biometrika* **104** 237–242. [MR3626474](#)
- [62] TZOURIO-MAZOYER, N., LANDEAU, B., PAPATHANASSIOU, D., CRIVELLO, F., ETARD, O., DELCROIX, N., MAZOYER, B. and JOLIOT, M. (2002). Automated Anatomical Labeling of Activations in SPM Using a Macroscopic Anatomical Parcellation of the MNI MRI Single-Subject Brain. *NeuroImage* **15** 273–289.
- [63] WAINWRIGHT, M. J. (2019). *High-Dimensional Statistics*. Cambridge University Press. [MR3967104](#)
- [64] WANG, B., WANG, G., WANG, X., CAO, R., XIANG, J., YAN, T., LI, H., YOSHIMURA, S., TOICHI, M. and ZHAO, S. (2021). Rich-club analysis in adults with ADHD connectomes reveals an abnormal structural core network. *Journal of Attention Disorders* **25** 1068–1079.
- [65] WANG, M., HU, Z., LIU, L., LI, H., QIAN, Q. and NIU, H. (2020). Disrupted functional brain connectivity networks in children with attention-deficit/hyperactivity disorder: evidence from resting-state functional near-infrared spectroscopy. *Neurophotonics* **7** 015012–015012.
- [66] WEI, Z. and LI, H. (2008). A hidden spatial-temporal Markov random field model for network-based analysis of time course gene expression data. *The Annals of Applied Statistics* **2** 408–429. [MR2415609](#)
- [67] ZAPATA, J., OH, S. Y. and PETERSEN, A. (2022). Partial separability and functional graphical models for multivariate Gaussian processes. *Biometrika* **109** 665–681. [MR4472841](#)

- [68] ZHAO, B., WANG, Y. S. and KOLAR, M. (2019). Direct Estimation of Differential Functional Graphical Models. In *Advances in Neural Information Processing Systems, NeurIPS*.
- [69] ZHAO, B., WANG, Y. S. and KOLAR, M. (2022). FuDGE: A Method to Estimate a Functional Differential Graph in a High-Dimensional Setting. *Journal of Machine Learning Research* **23** 1–82. [MR4576667](#)
- [70] ZHAO, B., ZHAI, P., WANG, Y. S. and KOLAR, M. (2024). Supplement to “High-dimensional Functional Graphical Model Structure Learning via Neighborhood Selection Approach”. *Electronic Journal of Statistics*, <https://doi.org/10.1214/24-EJS2219SUPP>.
- [71] ZHAO, P. and YU, B. (2006). On Model Selection Consistency of Lasso. *Journal of Machine Learning Research* **7** 2541–2563. [MR2274449](#)
- [72] ZHU, H., STRAWN, N. and DUNSON, D. B. (2016). Bayesian Graphical Models for Multivariate Functional Data. *Journal of Machine Learning Research* **17** 1–27. [MR3580357](#)