

Multi-relational Network Autoregression Model with Latent Group Structures

Yimeng Ren^{1,2}

YMREN@UST.HK

¹ *School of Data Science*

Fudan University

Shanghai, China

² *HKUST Business School*

The Hong Kong University of Science and Technology

Hong Kong SAR, China

Xuening Zhu*

XUENINGZHU@FUDAN.EDU.CN

School of Management

Fudan University

Shanghai, China

Ganggang Xu*

GANGXU@BUS.MIAMI.EDU

Department of Management Science

University of Miami

Coral Gables, FL 33146, USA

Yanyuan Ma

YANYUANMA@GMAIL.COM

Department of Statistics

The Pennsylvania State University

University Park, PA 16802, USA

Editor: Ji Zhu

Abstract

Multi-relational networks among entities are frequently observed in the era of big data. Quantifying the effects of multiple networks has attracted significant research interest recently. In this work, we model multiple network effects through an autoregressive framework for tensor-valued time series. To characterize the potential heterogeneity of the networks and handle the high dimensionality of the time series data simultaneously, we assume a separate group structure for entities in each network and estimate all group memberships in a data-driven fashion. Specifically, we propose a group tensor network autoregression (GTNAR) model, which assumes that within each network, entities in the same group share the same set of model parameters, and the parameters differ across networks. An iterative algorithm is developed to estimate the model parameters and the latent group memberships simultaneously. Theoretically, we show that the group-wise parameters and group memberships can be consistently estimated when the group numbers are correctly- or possibly over-specified. An information criterion for estimating the group number for each network is also provided to consistently select the group numbers. Lastly, we apply the GTNAR method to a Yelp dataset to illustrate its usefulness.

Keywords: Multi-relational networks, Latent group, Tensor-valued time series, Network autoregression.

* Corresponding authors.

1. Introduction

As the world becomes increasingly connected, studying network effects has become an important research topic across disciplines, including economics, finance, and many others. The primary focus of our study is to analyze the network effects inherent in high-dimensional time series observed over multiple networks. The existing literature has seen notable progress in the study of time series within a single network. For instance, Zhu et al. (2017) introduced a network autoregression model to investigate time series in large social networks. Armillotta and Fokianos (2023) extended the framework to the nonlinear network autoregression model, in which the multivariate observations could be both continuous and discrete. Chen et al. (2023) proposed a community network vector autoregression model for the high-dimensional time series. Fang et al. (2023) proposed a group network Hawkes process with a single network to model unit event occurrences. Chen et al. (2026) proposed a group network multivariate GARCH model that incorporates a latent group structure and an observed network adjacency matrix. Compared with conventional high-dimensional time series models (Walden and Serroukh, 2002; Leng and Tang, 2012; Zhou, 2014; Wang et al., 2019; Chen et al., 2020a; Chang et al., 2023; Wang and Tsay, 2023; Chen and Fan, 2023), the network autoregressive model distinguishes itself by providing enhanced parameter interpretability, offering deeper insights into the intricacies of network dynamics.

In real-world scenarios, entities within a population commonly establish connections across multiple networks, often referred to as multi-relational or multi-layer networks. Recently, there has been significant interest in investigating these networks collectively, as seen in works like Lei et al. (2020), Zhang et al. (2020), Jing et al. (2021), MacDonald et al. (2022), and Ma and Nandy (2023). While the majority of these studies focus on identifying community structures within multi-layer networks, there is also a considerable interest in quantifying the impacts of multi-relational network effects on various research objectives. For example, Emch et al. (2016) investigated the joint effects of spatial and social networks on disease transmission. Chen et al. (2017) proposed the utilization of network metrics from various social networks to predict the adoption of new products in marketing research. Besides, Corradini et al. (2021) investigated the influence of multi-dimensional social networks on negative reviews posted on Yelp. Although these models are valuable in empirical research, a critical gap remains in the availability of rigorous statistical models capable of providing valid statistical inferences about multiple network effects, and it is our intention to fill this gap.

The objective of our work is to investigate tensor-valued time series indexed across multi-relational networks. A straightforward way to analyze tensor-valued data is to stack it into vectors or to consider only one of its dimensions. However, this will destroy the intrinsic multidimensional structure and lack clear interpretations, and provide limited insights into network effects. This drawback has also been recognized in several recent studies on matrix- or tensor-valued time series (Zhou and Li, 2014; Wang et al., 2019; Chen et al., 2020a, 2021, 2022; Chen and Fan, 2023; Wang et al., 2024; Chen and Lam, 2024). Additionally, although stacking into vectors allows the implementation of techniques for high-dimensional VAR models, this still leads to a much larger number of parameters that need to be estimated, on the order of $\prod_{l=1}^q (N_l)^2$, where q is the dimension of tensor data, and N_l is the number of nodes in the l th dimension. Even with commonly used sparsity regularization,

the estimation variance of model parameters is likely to be much higher than that of the network VAR model when the data is generated by the latter (Zhu et al., 2023).

To address high dimensionality, several recent studies focus on factor models to capture low-rank structures (Zhou et al., 2023; Han et al., 2023; Wang et al., 2024). For example, Chen et al. (2022) introduced a general framework of factor models for tensor-valued time series. Han et al. (2024) propose a factor model and a high-order projection estimator to analyze high-dimensional dynamic tensor time series. Moreover, Wang et al. (2022) used the tensor decomposition technique to deal with the large transition matrices of the high-dimensional VAR model. Chen and Lam (2024) introduced a pre-averaging procedure, allowing for a spectrum of factor strengths, as well as the weak factors when both cross-sectional and serial correlations are present. Besides, Wang and Tsay (2023) proposed a general robust estimation procedure, which can address many high-dimensional VAR models with low-rank or sparse structure. However, existing works do not exploit network structure information across tensor dimensions and therefore do not provide direct statistical inference on these network effects.

On the other hand, it is remarkable that nodal heterogeneity widely exists in practice (Ke et al., 2015). The utilization of latent group structures to model heterogeneous data has a well-established history in panel data analysis. For instance, Ke et al. (2015) introduced a clustering algorithm in regression through data-driven segmentation to identify the groups, and Vogt and Linton (2017) shed light on the latent grouped structure in the non-parametric regression functions. Su et al. (2016) introduce a Classifier Lasso (C-Lasso) estimator for panel models. And more recently, Ando and Bai (2020) proposed a new estimation method for analyzing the quantile co-movement of large-scale financial time series data, which can identify latent group heterogeneity among the series. There have also been recent efforts to leverage latent group structures in modeling time series data within a single network, see, e.g., Zhu et al. (2023), Chen et al. (2023), Fang et al. (2023), Liu et al. (2024), and Chen et al. (2026). However, these works have predominantly focused on time series observed on a single network; hence, they cannot be directly used in modeling the tensor-valued time series data. Moreover, addressing the theoretical challenges involved in extending from a single network to multiple networks is non-trivial due to the interactions between different group structures, necessitating the development of new theoretical tools. To tackle this challenge, we divide the members of the network on each dimension into several latent subgroups. We assume that members within each subgroup share similar characteristics, carrying the same model coefficients.

The main contributions of our work can be summarized as follows. First, we introduce a highly interpretable network autoregression model for high-dimensional multivariate time series indexed by multiple networks, namely, the Group Tensor-valued Network Autoregression (GTNAR) model. Second, to account for network heterogeneity, we allow for separate group structures on related networks. Third, we establish estimation consistency for both model parameters and group memberships, not only when the numbers of groups are correct, but also when they are over-specified. Simultaneously estimating group memberships across multiple networks poses significant challenges compared to the existing literature, which typically considers at most one network (e.g., Fang et al., 2023). We develop new theoretical tools to address this challenge. Lastly, we develop a selection criterion that consistently selects the true group numbers and establishes asymptotic normality when the

group numbers are correctly specified. Our theoretical framework enables rigorous testing of multiple network effects, which are crucial across various research disciplines.

The remainder of this article is structured as follows. In Section 2, we first show a motivating example on the Yelp dataset to introduce a simple model form. Then, we introduce the general GTNAR model. Section 3 outlines the model estimation procedure and the selection method for the group numbers. Theoretical properties concerning parameter estimation, group membership estimation, and group number selection are discussed in Section 4. Two model extensions and corresponding theories are provided in Section 5. In Section 6, we present extensive simulation studies to illustrate the finite sample performance of GTNAR. Section 7 includes an application of GTNAR to the Yelp dataset. Finally, Section 8 provides concluding remarks. Additional technical proofs can be found in the Appendix.

2. Group Tensor-valued Network Autoregression

2.1 Notations and Tensor Algebra

Throughout the paper, we use the following notations. Denote $[n] = \{1, 2, \dots, n\}$ for an integer n . For a vector $\mathbf{v} = (v_j : j \in [p])^\top \in \mathbb{R}^p$, let $\|\mathbf{v}\| = (\sum_{j=1}^p v_j^2)^{1/2}$. For a matrix $\mathbf{M} = (m_{ij}) \in \mathbb{R}^{n_1 \times n_2}$, let \mathbf{M}_i be the i th row vector and \mathbf{M}_j as the j th column vector of \mathbf{M} . In addition, let $\mathbf{M}^{(\mathcal{C}, \cdot)} = (m_{ij} : i \in \mathcal{C}, j \in [n_2])$ and $\mathbf{M}^{(\cdot, \mathcal{C})} = (m_{ij} : i \in [n_1], j \in \mathcal{C})$, where \mathcal{C} is an index set. For a symmetric matrix \mathbf{M} , define $\lambda_{\min}(\mathbf{M})$ and $\lambda_{\max}(\mathbf{M})$ as the smallest and largest eigenvalues, respectively. For a tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_q}$, denote $\mathcal{X}^{(\mathcal{C}_1, \dots, \mathcal{C}_q)}$ as the corresponding subset of the tensor with each dimension selected by the index sets \mathcal{C}_l s. Denote $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{n_1 n_3 \times n_2 n_4}$ as the Kronecker product between matrices $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$ and $\mathbf{B} \in \mathbb{R}^{n_3 \times n_4}$. For a matrix series $\{\mathbf{A}_l : l \in [q]\}$, denote $\mathbf{B} \otimes_{k \neq l} \mathbf{A}_k = \mathbf{B} \otimes \mathbf{A}_1 \otimes \dots \otimes \mathbf{A}_{l-1} \otimes \mathbf{A}_{l+1} \otimes \dots \otimes \mathbf{A}_q$. For a vector, matrix, or tensor \mathbf{M} , let $\|\mathbf{M}\|_{\max}$ denote its largest absolute entry. For a set \mathcal{S} , denote $|\mathcal{S}|$ as the cardinal number of \mathcal{S} . We denote $\mathbf{1}_p \in \mathbb{R}^p$ as a p -dimensional vector with all elements equal to one, and \mathbf{I}_p is an identity matrix.

For a tensor $\mathcal{X} = (x_{i_1, \dots, i_q})_{i_1 \in [n_1], \dots, i_q \in [n_1]} \in \mathbb{R}^{n_1 \times \dots \times n_q}$, denote the vectorization as $\mathbf{x} = \text{vec}(\mathcal{X}) = (x_{1,1, \dots, 1}, x_{2,1, \dots, 1}, \dots, x_{n_1,1, \dots, 1}, x_{1,2,1, \dots, 1}, \dots, x_{n_1, n_2, \dots, n_q})^\top \in \mathbb{R}^{\prod_{l=1}^q n_l}$. Denote its mode- l matricization as $\mathcal{X}_{(l)} \in \mathbb{R}^{n_l \times \prod_{k \neq l} n_k}$, which is calculated by setting the l th tensor dimension as the matrix rows, and collapsing all other into its columns, for $l \in [q]$. For tensor $\mathcal{X} = (x_{i_1, \dots, i_q}) \in \mathbb{R}^{n_1 \times \dots \times n_q}$ and a matrix $\mathbf{M} = (m_{i,j}) \in \mathbb{R}^{q \times n_l}$, denote the mode- l multiplication $\mathcal{X} \times_l \mathbf{M}$ as a tensor in $\mathbb{R}^{n_1 \times \dots \times n_{l-1} \times q \times n_{l+1} \times \dots \times n_q}$, and its $(i_1, \dots, i_{l-1}, s, i_{l+1}, \dots, i_q)$ th element is calculated by $(\mathcal{X} \times_l \mathbf{M})_{i_1, \dots, i_{l-1}, s, i_{l+1}, \dots, i_q} = \sum_{i_l=1}^{n_l} x_{i_1, \dots, i_l, \dots, i_q} m_{s, i_l}$. Subsequently, we denote $\mathcal{X} \times_{l=1}^q \mathbf{M}_l \stackrel{\text{def}}{=} (\mathcal{X} \times_1 \mathbf{M}_1) \times_2 \dots \times_q \mathbf{M}_q$. For two tensors $\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_q}$ and $\mathcal{Y} \in \mathbb{R}^{n_1 \times \dots \times n_s}$ with $q \geq s$, their generalized inner product $\langle \mathcal{X}, \mathcal{Y} \rangle$ is a $(q-s)$ dimensional tensor, calculated by $\langle \mathcal{X}, \mathcal{Y} \rangle_{i_{s+1} \dots i_q} = \sum_{i_1=1}^{n_1} \dots \sum_{i_s=1}^{n_s} \mathcal{X}_{i_1 \dots i_s i_{s+1} \dots i_q} \mathcal{Y}_{i_1 \dots i_s}$. We use $\mathcal{A} \odot \mathcal{Y}$ to denote the element-wise product between two tensors (\mathcal{A} and \mathcal{Y}) with the same dimensions. For a column vector $\mathbf{x}_l \in \mathbb{R}^{n_l}$, $\mathbf{x}_l \circ_{k \neq l} \mathbf{1}_{n_k} = \mathbf{1}_{n_1} \circ \dots \circ \mathbf{x}_l \circ \dots \circ \mathbf{1}_{n_q}$ is a q dimensional tensor $\mathcal{T} \in \mathbb{R}^{n_1 \times \dots \times n_q}$ with its $(i_1, \dots, i_l, \dots, i_q)$ th element being x_{i_l} .

2.2 Motivating Example and A Simple Model

The objective of our work is to investigate time series indexed across multi-relational networks. To offer a more lucid representation of GTNAR, we demonstrate it through a motivating dataset collected from Yelp (<https://www.yelp.com/>). Yelp serves as a prominent review platform for various businesses, including restaurants, local retailers, entertainment establishments, and more. It also functions as a social platform where users can share information and their personal experiences. The dataset covers the period from 2010 to 2018 and is collected from five North American cities (i.e., Charlotte, Las Vegas, Phoenix, Scottsdale and Toronto), and comprises four main categories of information: user data (e.g., user registration time on Yelp), user-friend relationships, business information (including spatial location), and user reviews of businesses. For example, Figure 1 illustrates a user’s review of a restaurant named “Esther’s Kitchen” in Las Vegas. In this case, the user gave the restaurant a five-star rating, and their review received 18 tags from other users, including 8 “useful”, 3 “funny”, and 7 “cool” tags. Overall, the restaurant has accumulated 1611 reviews.

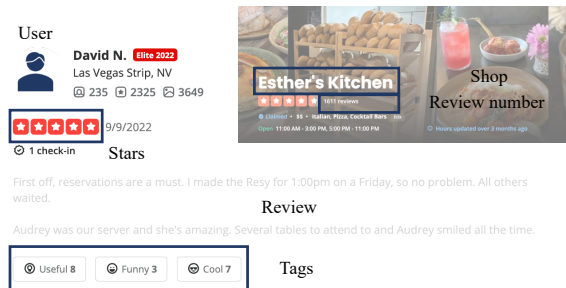


Figure 1: A review snapshot of the shop “Esther’s Kitchen”. It contains the user information, shop statistics, review text, and the tags assigned to this review.

Within the Toronto segment of the dataset, we have records for $N_1 = 462$ users who have provided reviews for restaurants in $N_2 = 56$ distinct locations, spanning $T = 36$ quarters. Our primary focus in this analysis concentrates on the variable denoted as $Y_{ij,t}$, representing the $\log(1+x)$ -transformed number of reviews contributed by user i to restaurants in district j during the t th quarter. This variable forms a time series indexed by both the user ID (i) and the district ID (j). The first challenge we encounter when analyzing this dataset is that neither users nor districts can be considered isolated units. As a result, it becomes imperative to model the $Y_{ij,t}$ ’s in a collective manner. Specifically, users form a social network, while districts establish a spatial network that fosters substantial interactions among their respective network members. These interactions, in turn, significantly influence the outcome variable $Y_{ij,t}$ when considered jointly. For example, Tiwari and Richards (2016) found that peer social networks play a highly effective role in influencing restaurant preferences within social circles of friends. The social network analysis conducted on Yelp data by Fe (2023) suggests that social network friends are 64% more likely to visit the same restaurant when compared to non-friends. Sun and Paule (2017) discovered significant

spatial effects on ratings across various categories of Yelp venues, and Gan et al. (2021) investigated the spatial network effects on the tourism economy. However, these studies are primarily empirical and lack a strong statistical foundation. Furthermore, they tend to concentrate solely on the impact of a single network, rather than considering multiple network effects jointly. As a motivating example, Figure 2(a) depicts the social network of Toronto users with at least two friends, highlighting the observation that connected friends often have similar comment volumes. Meanwhile, Figure 2(b) presents the spatial network of Toronto, indicating that neighboring districts, such as zones 1 and 2, tend to exhibit similar comment volumes. Both network effects contribute jointly to the outcome variable. Therefore, the first challenge we intend to address is how to construct a multivariate time series model that can rigorously quantify the impact of multiple network effects for data similar to the Yelp review dataset.

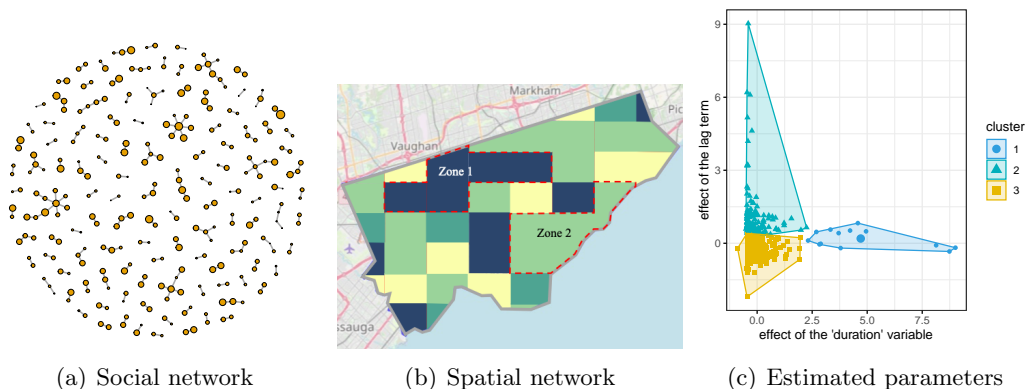


Figure 2: (a) The social network of Toronto users with degrees greater than one. Node sizes reflect the logarithm of the number of comments made in 2008. (b) The spatial network of Toronto districts during the last quarter of 2018. Colors represent four quantile intervals of the logarithm of the number of comments made in the 4th quarter of 2008. Darker colors indicate districts with higher comment activity. (c) Clustering results based on estimated regression coefficients, categorizing users into three distinct groups. Different colors and shapes visually represent each group.

The second challenge we encounter in the Yelp review dataset relates to the heterogeneity among members in both social and spatial networks. Previous research in network analysis (Newman et al., 2006) highlighted that users tend to naturally form clusters based on similar behaviors, a phenomenon consistent with Yelp’s ecosystem, and users display distinct interaction patterns (e.g., frequent reviewers versus casual browsers) (Fe, 2023; Sun and Paule, 2017). There is also empirical research (Wedel and Kamakura, 2000; Fiebig et al., 2010), demonstrating that segmenting users enhances recommendation accuracy on platforms with varied preferences, which is particularly applicable to Yelp’s diverse range of restaurants and businesses (Jagabathula et al., 2018). Similarly, restaurants located in different spatial regions, such as the central business district or other areas, may experience different levels of popularity, leading to varying degrees of spatial spillover effects (see, e.g., Koschinsky, 2009). Recent advances (e.g., Bonhomme and Manresa (2015); Su et al. (2016);

Liu et al. (2020)) in high-dimensional time series modeling demonstrated that subgrouping mitigates bias in parameter estimation and enhances estimation efficiency given limited sample sizes, which is also a critical feature for the Yelp data. As a motivating example, we conduct a preliminary regression analysis with aggregated users' comment volumes (log-transformed) as the response variable, denoted as $Y_{it} = \sum_j Y_{ij,t}$. Two covariates are included: the lagged response ($Y_{i(t-1)}$) and the user's duration since registration (i.e., the number of months since joining Yelp as of the $(t-1)$ th quarter). The estimated regression coefficients obtained from all users are clustered into three groups using the k -means algorithm, and these clusters are visualized in Figure 2(c). An evident heterogeneous pattern is observed among the users' coefficients. To tackle this challenge, we adopt the approach proposed in Zhu et al. (2023), which involves dividing the members of both social and spatial networks into several subgroups. We assume that members of each subgroup share similar characteristics.

Let $\mathbf{A}^{(1)} = (a_{ij}^{(1)}) \in \mathbb{R}^{N_1 \times N_1}$ and $\mathbf{A}^{(2)} = (a_{ij}^{(2)}) \in \mathbb{R}^{N_2 \times N_2}$ represent the observed exogenous adjacency matrices characterizing the social and spatial networks, respectively. Specifically, $a_{ij}^{(1)} = 1$ implies that the i th user follows the j th user, while $a_{ij}^{(1)} = 0$ otherwise. Similarly, $a_{ij}^{(2)} = 1$ indicates that the i th district is a spatial neighbor of the j th district, while $a_{ij}^{(2)} = 0$ otherwise. We follow the convention by setting $a_{ii}^{(1)} = 0$ and $a_{jj}^{(2)} = 0$ for $1 \leq i \leq N_1$ and $1 \leq j \leq N_2$. We assume the existence of G_1 groups in the social network and G_2 groups in the spatial network. The group membership for the i th node in the social network is denoted as $g_i^{(1)}$ ($1 \leq g_i^{(1)} \leq G_1$), and the group membership for the j th node in the spatial network is denoted as $g_j^{(2)}$ ($1 \leq g_j^{(2)} \leq G_2$). We propose the following model with a two-way group structure:

$$\begin{aligned}
 Y_{ij,t} = & \underbrace{\lambda_{g_i^{(1)}}^{(1)} \sum_{k=1}^{N_1} \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,(t-1)}}_{\text{Social Network main effect}} + \underbrace{\lambda_{g_j^{(2)}}^{(2)} \sum_{k=1}^{N_2} \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,(t-1)}}_{\text{Spatial Network main effect}} \\
 & + \underbrace{\alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)}}_{\text{Self-momentum}} + \underbrace{\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)}}_{\text{Covariate effects}} + \varepsilon_{ij,t},
 \end{aligned} \tag{1}$$

where $n_{1i} = \sum_{k=1}^{N_1} a_{ik}^{(1)}$, $n_{2j} = \sum_{k=1}^{N_2} a_{kj}^{(2)}$, $\mathbf{x}_{it}^{(1)} \in \mathbb{R}^{p_1}$ and $\mathbf{x}_{jt}^{(2)} \in \mathbb{R}^{p_2}$ are exogenous covariate vectors of finite dimensions associated with the i th user and j th district, respectively. Besides, $\varepsilon_{ij,t}$ represents independent and identically distributed (i.i.d.) white noise with $E(\varepsilon_{ij,t}) = 0$ and its variance $\text{var}(\varepsilon_{ij,t}) = \sigma^2$. For identifiability, it is required that $\sum_{g^{(1)}=1}^{G_1} \zeta_{g^{(1)},1}^{(1)} = 0$ when both intercepts are included in $\mathbf{x}_{it}^{(1)}$ and $\mathbf{x}_{it}^{(2)}$, where $\zeta_{g_i^{(1)},1}^{(1)}$ (the first element of $\boldsymbol{\zeta}_{g_i^{(1)}}^{(1)}$) represents the intercept for $\mathbf{x}_{it}^{(1)}$.

The first term in (1), i.e., $\sum_{k=1}^{N_1} (a_{ik}^{(1)}/n_{1i}) Y_{kj,(t-1)}$, represents the average number of reviews (log(1 + x)-transformed) by user i 's following friends on the restaurants in district j in the previous quarter. Consequently, $\lambda_{g_i^{(1)}}^{(1)}$ quantifies the influence of following friends on user i 's attitude towards district j and encapsulates a social network main effect. On the other hand, the second term in (1), i.e., $\sum_{k=1}^{N_2} (a_{kj}^{(2)}/n_{2j}) Y_{ik,(t-1)}$, calculates the average

log number of reviews by the i th user on districts that are connected with the j th district in the previous quarter. Thus, $\lambda_{g_j^{(2)}}^{(2)}$ signifies how the user's review count towards district j is influenced by his/her reviews towards the districts connected with district j , and can be interpreted as the spatial network main effect. Additionally, $\alpha_{g_i^{(1)}g_j^{(2)}}$ represents the self-driven time effect for the (i, j) th time series, quantifying the momentum effect of the review activity by the i th user towards the j th district in the previous quarter. A higher value of $\alpha_{g_i^{(1)}g_j^{(2)}}$ suggests a greater level of loyalty of user i to district j . Finally, $\zeta_{g_i^{(1)}}^{(1)} \in \mathbb{R}^{p_1}$ and $\zeta_{g_j^{(2)}}^{(2)} \in \mathbb{R}^{p_2}$ are external covariate effects at the user and district levels, enhancing the model's capacity to account for user and district heterogeneity in the data. By evaluating the significance of $\lambda_{g^{(1)}}^{(1)}$'s and $\lambda_{g^{(2)}}^{(2)}$'s, one can examine the presence of social and spatial network main effects while controlling for other factors in Model (1).

Our framework for multiple networks further extends the two-network model (1) to multi-relational networks given by Model (2) in Section 2.3 below, which we refer to as the GTNAR model.

2.3 General Model in Tensor Form

Suppose there are q observed exogenous networks characterized by adjacency matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(q)}$, with the l th matrix $\mathbf{A}^{(l)} = (a_{ij}^{(l)}) \in \mathbb{R}^{N_l \times N_l}$, where N_l is the number of nodes in the l th network, and q is finite. Specifically, $a_{ij}^{(l)} = 1$ implies that the i th node is connected with the j th node in the l th network, while $a_{ij}^{(l)} = 0$ otherwise. We follow the convention by setting $a_{i_l i_l}^{(l)} = 0$ for $1 \leq i_l \leq N_l$ and $1 \leq l \leq q$. We assume that there exists G_l groups in the l th network, and the group membership for the i th node in this network is denoted as $g_i^{(l)}$ ($1 \leq g_i^{(l)} \leq G_l$). As a result, each network has its own latent group structure among the N_l network nodes, which is denoted by $\mathcal{G}_l = (g_i^{(l)} : 1 \leq i \leq N_l)$, $1 \leq l \leq q$. Under such a model framework, the response variables of interest constitute a tensor-valued time series, denoted by $\mathcal{Y}_t = (Y_{i_1 i_2 \dots i_q, t}) \in \mathbb{R}^{N_1 \times N_2 \times \dots \times N_q}$, with the following model structure

$$\begin{aligned}
 Y_{i_1 i_2 \dots i_q, t} = & \sum_{l=1}^q \lambda_{g_{i_l}^{(l)}}^{(l)} \underbrace{\sum_{k=1}^{N_l} \frac{a_{i_l k}^{(l)}}{n_{l i_l}} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)}}_{\text{The } l\text{th Network main effect}} + \underbrace{\alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} Y_{i_1 i_2 \dots i_q, (t-1)}}_{\text{Self-momentum}} \\
 & + \sum_{l=1}^q \underbrace{\mathbf{x}_{i_l t}^{(l)\top} \zeta_{g_{i_l}^{(l)}}^{(l)}}_{\text{The } l\text{th covariate effects}} + \varepsilon_{i_1 i_2 \dots i_q, t}, \quad (2)
 \end{aligned}$$

where $n_{l i_l} = \sum_{k=1}^{N_l} a_{i_l k}^{(l)}$, $\mathbf{x}_{i_l t}^{(l)} \in \mathbb{R}^{p_l}$ represents exogenous covariates associated with the i_l th member in the l th network, and $\varepsilon_{i_1 i_2 \dots i_q, t}$ denotes the white noise with $\text{var}(\varepsilon_{i_1 \dots i_q, t}) = \sigma^2$. In this work, we treat the number of layers (q) and the covariates dimensions (p_l , $l \in [q]$) as fixed constants. For indentifiability, assume that $\sum_{g^{(l)=1}}^{G_l} \zeta_{g^{(l),1}}^{(1)} = 0$ ($l \in [q-1]$). Define $\mathcal{G}_l = (g_{i_l}^{(l)} : 1 \leq i_l \leq N_l)^\top \in \mathbb{R}^{N_l}$ and denote $\mathcal{G} = \{\mathcal{G}_l : 1 \leq l \leq q\}$. Let $\mathcal{R}_g^{(l)} = \{i_l : g_{i_l}^{(l)} = g\}$

and further denote $N_{lg} = |\mathcal{R}_g^{(l)}|$. Denote $g^{-(l)} = (g^{(1)}, \dots, g^{(l-1)}, g^{(l+1)}, \dots, g^{(q)})$ as the group membership indices excluding the l th network, and denote the group index $\mathcal{I}_{g^{(1)}, \dots, g^{(q)}}$ as a function of $(g^{(1)}, \dots, g^{(q)})$, which takes values from $\{1, \dots, \prod_l G_l\}$.

Define $\mathbf{W}^{(l)} = (a_{ij}^{(l)}/n_{li})$ as the l th row-normalized adjacency matrix of $\mathbf{A}^{(l)}$. Then the GTNAR model (2) can be expressed in a tensor form as follows,

$$\mathcal{Y}_t = \sum_{l=1}^q (\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)}) \times_l \mathbf{L}^{(l)} + \mathcal{A} \odot \mathcal{Y}_{t-1} + \sum_{l=1}^q \boldsymbol{\beta}_{X_l, t}^{(l)} \circ_{k \neq l} \mathbf{1}_{N_k} + \mathcal{E}_t, \quad (3)$$

where $\mathbf{L}^{(l)} = \text{diag}(\lambda_{g_{i_l}}^{(l)} : 1 \leq i_l \leq N_l) \in \mathbb{R}^{N_l \times N_l}$, $\mathcal{A} = (\alpha_{g_{i_1} \dots g_{i_q}} : 1 \leq i_l \leq N_l) \in \mathbb{R}^{N_1 \times \dots \times N_q}$, $\boldsymbol{\beta}_{X_l, t}^{(l)} = (\mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}}^{(l)} : 1 \leq i_l \leq N_l)^\top \in \mathbb{R}^{N_l}$ and $\mathcal{E}_t = (\varepsilon_{i_1, \dots, i_q, t}) \in \mathbb{R}^{N_1 \times \dots \times N_q}$.

Remark 1. (Matrix Autoregression When $q = 2$) When $q = 2$, the model in (3) can be simplified as the matrix-valued autoregression model in (4) below,

$$\mathbf{Y}_t = (\mathbf{L}^{(1)} \mathbf{W}^{(1)}) \mathbf{Y}_{t-1} + \mathbf{Y}_{t-1} (\mathbf{W}^{(2)} \mathbf{L}^{(2)}) + \mathbf{A} \circ \mathbf{Y}_{t-1} + \boldsymbol{\beta}_{X_1, t} \mathbf{1}_{N_2}^\top + \mathbf{1}_{N_1} \boldsymbol{\beta}_{X_2, t}^\top + \mathbf{E}_t, \quad (4)$$

where $\mathbf{Y}_t \in \mathbb{R}^{N_1 \times N_2}$ is the response matrix, $\mathbf{L}^{(l)} = \text{diag}(\lambda_{g_{i_l}}^{(l)} : 1 \leq i_l \leq N_l) \in \mathbb{R}^{N_l \times N_l}$ for $l = 1, 2$, $\mathbf{A} = (\alpha_{g_{i_1} g_{i_2}} : 1 \leq i_1 \leq N_1, 1 \leq i_2 \leq N_2) \in \mathbb{R}^{N_1 \times N_2}$, and $\mathbf{E}_t \in \mathbb{R}^{N_1 \times N_2}$ is the random noise matrix.

Remark 2. We remark that GTNAR is not limited to the motivating example and can be utilized in a wide range of applications. Specifically, GTNAR is designed for tensor-valued time series data, provided we can collect network relationships among the units. The network under consideration is not restricted to spatial networks, as we use in the motivating example. For instance, in financial data analysis, we can use the common shareholding relationship to examine the spillover effects of systematic risk (Feng et al., 2023). In the international trading market, regions can construct two adjacent networks from the export and import aspects, while the trading products categories form a similarity network (Alves et al., 2019). In the contextual recommendation system, the social network relationships among the users form a network on the first dimension, the similarities based on the item characteristics form the second one, and the proximities of different types of users' activities mode form the third network (Adomavicius et al., 2005; Wu et al., 2016). Other examples include supply chain linkages (Serpa and Krishnan, 2018), common interests in customer purchase (Carroni et al., 2020), and others, according to the application scenarios.

Remark 3. Applications of tensor-valued time series are widespread, spanning numerous fields. For instance, a contextual recommender system naturally generates a 3-dimensional tensor of users, items, and contexts, as it recommends items to users under varying contexts (such as the promotion types and the users' locations) over a period of time; each dimension can be further enriched by its corresponding network, such as a user social network, an item similarity network, and a contextual proximity network. Similarly, international trade data, capturing the export volumes of different product categories between regions over a period of time, forms another 3-dimensional tensor-valued time series, where geographical

and product-similarity networks can be constructed on each dimension. In the domain of neuroimaging, functional Magnetic Resonance Imaging (fMRI) produces a series of 3D brain scans, creating a 3-dimensional tensor-valued time series where networks can be defined by the adjacency of brain regions. A fourth example is found in online live-streaming, where user activities (such as comments and shares) directed at various streamers and different topics form a 4-dimensional tensor. This type of tensor recorded along a period of time naturally forms a 4-dimensional tensor-valued time series. The practical and urgent need to analyze this type of complex, high-dimensional data, which is ubiquitous in modern socio-economic and scientific contexts, serves as the primary applied motivation for the GTNAR model.

Remark 4. (Parameter Dimension Reduction) We remark that we reduce the parameter dimension by embedding the observed network weighting matrices $\mathbf{W}^{(k)}$ and a group structure. Take the case with $q = 2$ for example. In general, one needs to estimate $O(N_1^2 + N_2^2)$ parameters in model (4). In contrast, with the imposed model structure, the number of estimated parameters reduces to $O(G_1(p_1 + 1) + G_2(p_2 + 1) + G_1G_2)$ when the group memberships are given, which is greatly reduced due to incorporating the weighting matrices $\mathbf{W}^{(l)}$ and memberships \mathcal{G}_1 and \mathcal{G}_2 .

2.4 Comparisons with Existing Literature

There are numerous works that are closely related to GTNAR. The most pertinent works encompass high-dimensional vector autoregression models (Davis et al., 2016; Wang et al., 2022; Miao et al., 2023), tensor-valued time series models (Hoff, 2015; Wang et al., 2019; Chen et al., 2021, 2022; Chen and Fan, 2023; Wang et al., 2024), grouped panel data models (Ke et al., 2015; Ando and Bai, 2016; Su et al., 2016), and the network autoregressive models with heterogeneous effects (Zhu and Pan, 2020; Zhu et al., 2023). In this subsection, we illustrate how the GTNAR model differs from the above existing models.

2.4.1 HIGH-DIMENSIONAL VECTOR AUTOREGRESSION MODELS

Recently, there have been studies focusing on the high-dimensional vector autoregression (VAR) models (Davis et al., 2016; Wang et al., 2022; Miao et al., 2023). By stacking elements of the tensor $\mathcal{Y}_t = (Y_{i_1 \dots, i_q, t}) \in \mathbb{R}^{N_1 \times \dots \times N_q}$ into a vector, GTNAR model (3) can be re-written as a VAR model with dimension $N' = \prod_l N_l$,

$$\mathbf{y}_t = \left[\sum_{l=1}^q \left((\mathbf{L}^{(l)} \mathbf{W}^{(l)}) \otimes_{k \neq l} \mathbf{I}_{N_k} \right) + \text{diag}(\mathbf{a}) \right] \mathbf{y}_{t-1} + \sum_{l=1}^q \left(\boldsymbol{\beta}_{X_l, t}^{(l)} \otimes_{k \neq l} \mathbf{1}_{N_k} \right) + \boldsymbol{\varepsilon}_t,$$

where $\mathbf{a} = \text{vec}(\mathcal{A}) \in \mathbb{R}^{N_1 \times \dots \times N_q}$ and $\boldsymbol{\varepsilon}_t = \text{vec}(\boldsymbol{\varepsilon}_t)$. This allows for the application of existing techniques of high-dimensional VAR models. While this formulation imposes fewer assumptions on the model coefficients, it leads to a much larger number of parameters to be estimated. Specifically, the number of parameters is $O(\prod_l N_l^2)$, in comparison to $O(\sum_l G_l(p_l + 1) + \prod_l G_l)$ (when fixing the group memberships) for the GTNAR model. Even with commonly used sparsity regularization, the estimation variance is likely to inflate compared to the network models under certain scenarios (Zhu et al., 2023), which is also investigated in our numerical studies. In addition, the “stacking” operation will destroy the

intrinsic data structure in the original tensor form, hence the coefficients estimated in the general VAR model lack clear interpretations and provide limited insights into the network effects. These drawbacks are also recognized in existing literatures (Chen et al., 2021; Wang et al., 2024).

2.4.2 TENSOR-VALUED TIME SERIES MODELS

This line of research contains two main categories, the tensor autoregressive framework (Ding and Dennis Cook, 2018; Chen et al., 2021; Wang et al., 2024) category, and the tensor factor model (Hoff, 2015; Wang et al., 2019; Chen et al., 2020b, 2022; Chen and Fan, 2023) category. Specifically, GTNAR should be classified into the first category. For ease of understanding, we compare our modeling approach with existing approaches when $q = 2$, which reduces to the matrix case in (4).

First, the tensor autoregressive (TAR) models characterize the dynamics of \mathbf{Y}_t using its lagged information \mathbf{Y}_{t-1} . For example, the TAR model proposed by Wang et al. (2024) takes the following form, i.e.,

$$\mathcal{Y}_t = \langle \mathcal{A}, \mathcal{Y}_{t-1} \rangle + \mathcal{E}_t, \quad (5)$$

where $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_q \times n_1 \times \dots \times n_q}$ is a $2q$ dimensional tensor, which has Tucker ranks (r_1, \dots, r_q) with $r_l = \text{rank}(\mathcal{A}_{(l)})$. They assume that the tensor \mathcal{A} takes a low rank structure as $\mathcal{A} = \mathcal{C} \times_{l=1}^{2q} \mathbf{U}_l$, $\mathcal{C} \in \mathbb{R}^{r_1 \times \dots \times r_{2q}}$ is the core low-rank tensor and $\mathbf{U}_l \in \mathbb{R}^{n_l \times r_l}$ for $l \in [2q]$. When $q = 2$, the model (5) can be rewritten as

$$\mathbf{Y}_t = \mathbf{U}_3(\mathcal{C}\mathbf{U}_1^\top \mathbf{Y}_{t-1} \mathbf{U}_2) \mathbf{U}_4^\top + \mathbf{E}_t. \quad (6)$$

In model (6), the number of parameters to be estimated is $O\{\prod_{l=1}^4 r_l + 2(N_1 + N_2)\}$. Compared to the existing TAR modeling approach, the GTNAR model (4) has two major differences. First, we embed the network weighting matrices $\mathbf{W}^{(l)}$ into the autoregression matrices to characterize the network dependence structure. To further reduce the number of estimated parameters, we consider a group structure with autoregression coefficients in $\mathbf{L}^{(1)}$ and $\mathbf{L}^{(2)}$, respectively. As we commented in Remark 4 for $q = 2$, GTNAR contains $O(G_1(p_1 + 1) + G_2(p_2 + 1) + G_1 G_2)$ parameters in the network autoregression coefficients when the group memberships are fixed. In addition, we provide an interpretable modeling framework with the group-specific autoregression coefficients. Furthermore, we would like to remark that the autoregression matrices in model (4) cannot be expressed with a low rank form unless the network matrices $\mathbf{W}^{(l)}$ take a low rank structure. Second, in addition to the network autoregression terms, we use $\mathbf{A} \circ \mathbf{Y}_{t-1}$ to capture the self-momentum effect, which cannot be characterized by (5) with the low-rank assumption. One can require \mathbf{A} to have a specific low rank structure $\mathbf{A} = \mathbf{M}^{(1)} \mathbf{G} \mathbf{M}^{(2)\top}$, where $\mathbf{M}^{(l)} = (\mathbf{e}_{g_i^{(l)}} : 1 \leq i \leq N_l)^\top \in \mathbb{R}^{N_l \times G_l}$ is a membership matrix for the l -th layer, and $\mathbf{G} = (\alpha_{gg'} : g \in [G_1], g' \in [G_2]) \in \mathbb{R}^{G_1 \times G_2}$. This relates GTNAR to the traditional low-rank models (Chen et al., 2021; Wang et al., 2024). We remark that the self-momentum effect contains cross-layer parameters in \mathbf{G} , which poses challenges in our theoretical analysis since the memberships from each layer cannot be analyzed separately.

The second category of tensor-valued time series models includes the tensor factor (TF) models considered by Wang et al. (2019), Chen et al. (2022), and Chen and Fan (2023).

The TF model assumes that the dynamics of \mathcal{Y}_t can be driven by a low-dimensional latent dynamic tensor factor \mathcal{F}_t . For example, the TF model proposed by Chen et al. (2022) takes the form $\mathcal{Y}_t = \mathcal{F}_t(\prod_k \times_k) \mathbf{A}_k + \mathcal{E}_t$. We refer to Hoff (2015) and Chen et al. (2020b) for other TF models in various forms. Although the TF model can be used to directly conduct a dimension reduction for \mathcal{Y}_t , it cannot quantify the historical self-momentum effect for \mathcal{Y}_t , and it also lacks a sound interpretation regarding the latent group structure.

We also compare the numerical performance of GTNAR with that of the multilinear tensor model in Hoff (2015) to illustrate the advantages of our proposed model.

2.4.3 GROUPED PANEL DATA MODELS

Utilizing latent group structures to model heterogeneous data has a well-established history in statistical models (Ke et al., 2015; Bester and Hansen, 2016; Liu et al., 2020). For instance, Bonhomme and Manresa (2015) and Bester and Hansen (2016) introduced grouped linear panel models with time-varying fixed effects and individual fixed effects, respectively. Su et al. (2016) introduced a Classifier Lasso (C-Lasso) estimator for panel models, and more recently, Liu et al. (2020) explored estimation and inference in the presence of possible over-specification of the group number. Specifically, Su et al. (2016) considered the grouped linear penal model $Y_{it} = \beta_i^\top \mathbf{x}_{it} + \mu_i + \varepsilon_{it}$, where the parameter β_i takes values from G_0 group centers $\{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_{G_0}\}$. However, the model format designed for vector-valued responses limits its applicability to complex tensor-valued data, thereby lacking generality. Compared to these studies, the GTNAR model further shows the additional advantage in quantifying the multiple heterogeneous network effects by incorporating different types of networks in tensor dimensions.

2.4.4 NETWORK AUTOREGRESSION MODEL WITH GROUP HETEROGENEITY

There have also been recent efforts to leverage latent group structures in modeling time series data within a single network, as demonstrated by Zhu and Pan (2020); Zhu et al. (2023); Chen et al. (2023); and etc. However, these works have predominantly focused on time series observed on a single network, and to the best of our knowledge, our work is the first to tackle time series indexed by multiple networks with distinct group structures. Specifically, Zhu et al. (2023) considered an autoregression model for a vector formed response \mathbf{y}_t as

$$\mathbf{y}_t = \mathbf{B}\mathbf{y}_{t-1} + \boldsymbol{\mu}_z + \boldsymbol{\varepsilon}_t, \quad (7)$$

where $\mathbf{B} \in \mathbb{R}^{N \times N}$ with its (i, j) th element being $b_{ij} = w_{ij}\beta_{g_i g_j}$ with $w_{ij} = a_{ij}/n_i$ for $i \neq j$ and $b_{ii} = \nu_{g_i}$ for $i = 1, \dots, N$. In (7), $\boldsymbol{\mu}_z = (\mathbf{x}_1^\top \boldsymbol{\zeta}_{g_1}, \dots, \mathbf{x}_N^\top \boldsymbol{\zeta}_{g_N})^\top$ represents the covariates term. The key difference between these two approaches lies in how they model network interactions. The goal of Zhu et al. (2023) is to capture interactions within a single network as flexibly as possible, where the effect between nodes i and j , namely, $w_{ij}\beta_{g_i g_j}$, depends jointly on their group memberships g_i and g_j . This makes the consistent estimation of group memberships substantially more challenging than in other existing group panel data models (Su et al., 2016; Liu et al., 2020), where group memberships do not interact in the parameter structure. To address this difficulty, Zhu et al. (2023) proposed a membership refinement procedure based on brute-force enumeration. While a straightforward extension of this approach to the multiple network tensor-valued data setting studied might be possible, we

adopt a completely different within-network effect model in (4), where in each network, the effect between node i and j takes the form $w_{ij}^{(l)}\lambda_{g_i^{(l)}}$, depending only on the membership of node i (but not j) within the l th network, $l = 1, \dots, q$. This approach has the advantage of avoiding the refinement procedure of Zhu et al. (2023), while has its own complexity of handling multiple networks simultaneously. Thus, the theoretical challenges in Zhu et al. (2023) and in our work stem from different sources, and we employ very different strategies to resolve them. The GTNAR model is a new modeling approach that leverages the general grouping idea to address challenges unique to the multi-layer networks setting.

3. Model Estimation

We first introduce the necessary notations. For the network corresponding to the l th tensor dimension, there exist G_l groups for the network nodes, denoted by $g^{(l)} \in [G_l]$. For each group, the group-level parameters are denoted by $\boldsymbol{\theta}_{g^{(l)}}^{(l)} = (\lambda_{g^{(l)}}^{(l)}, \boldsymbol{\zeta}_{g^{(l)}}^{(l)\top})^\top \in \mathbb{R}^{p_l+1}$. The collection of layer-specific parameters is then denoted as $\boldsymbol{\theta}^{(l)} = (\boldsymbol{\theta}_{g^{(l)}}^{(l)} : g^{(l)} \in [G_l]) \in \mathbb{R}^{G_l(p_l+1)}$, leading to the collection of $\boldsymbol{\xi} = (\boldsymbol{\theta}^{(1)\top}, \dots, \boldsymbol{\theta}^{(q)\top}, \text{vec}(\boldsymbol{\alpha})^\top)^\top \in \mathbb{R}^{\sum_l G_l(p_l+1) + \prod_l G_l}$ for all q tensor dimensions. Finally, the self-momentum parameter is denoted as $\boldsymbol{\alpha} = (\alpha_{g^{(1)}\dots g^{(q)}})_{g^{(l)} \in [G_l]} \in \mathbb{R}^{G_1 \times \dots \times G_q}$. In the following, we discuss estimation for the GTNAR model, which includes the group parameters $\boldsymbol{\xi}$ and memberships \mathcal{G} . We utilize an iterative algorithm to update $\boldsymbol{\xi}$ and \mathcal{G} . We first discuss the estimation of $\boldsymbol{\xi}$ when \mathcal{G} is given, and then we introduce the iterative algorithm for joint estimation of $\{\boldsymbol{\xi}, \mathcal{G}\}$. For easy understanding, we first introduce the estimation procedure for the matrix model in (4) when $q = 2$. Then we present a general estimation procedure for a general q .

3.1 Estimation Procedure for Model (4) when $q = 2$

We first introduce the estimation of model (4) with two networks for clear illustration. We aim to minimize the following objective function

$$Q(\boldsymbol{\xi}, \mathcal{G}) = \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} \sum_{t=1}^T \left(Y_{ij,t} - \lambda_{g_i^{(1)}}^{(1)} \sum_{k=1}^{N_1} \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,(t-1)} - \lambda_{g_j^{(2)}}^{(2)} \sum_{k=1}^{N_2} \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,(t-1)} - \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)} - \mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} - \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} \right)^2. \quad (8)$$

To minimize (8), we solve

$$\frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \boldsymbol{\theta}_{g^{(1)}}^{(1)}} = \mathbf{0}, \quad \frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \boldsymbol{\theta}_{g^{(2)}}^{(2)}} = \mathbf{0}, \quad \frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \alpha_{g^{(1)} g^{(2)}}} = 0.$$

By solving the equations above, we obtain the estimator when $q = 2$ as $\widehat{\boldsymbol{\theta}} = \mathbf{M}^{-1}\boldsymbol{\beta}$, where

$$\mathbf{M} = \begin{pmatrix} \mathbf{M}^{(1)} & \mathbf{M}^{(12)} & \mathbf{M}^{(1\alpha)} \\ \mathbf{M}^{(12)\top} & \mathbf{M}^{(2)} & \mathbf{M}^{(2\alpha)} \\ \mathbf{M}^{(1\alpha)\top} & \mathbf{M}^{(2\alpha)\top} & \mathbf{M}^\alpha \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \mathbf{b}^\alpha \end{pmatrix},$$

and the specific expressions are given in Appendix B.

3.2 Estimation Procedure for General Model with General q

We now discuss the estimation of the GTNAR model (3) in general. To simultaneously estimate the model parameters and the group memberships, we aim to minimize the following least squares objective function:

$$Q(\boldsymbol{\xi}, \mathcal{G}) = \sum_{i_1=1}^{N_1} \cdots \sum_{i_q=1}^{N_q} \sum_{t=1}^T \left(Y_{i_1 i_2 \dots i_q, t} - \sum_{l=1}^q \lambda_{g_{i_l}}^{(l)} \sum_{k=1}^{N_l} \frac{a_{i_l k}^{(l)}}{n_{l i_l}} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)} \right. \\ \left. - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} Y_{i_1 i_2 \dots i_q, (t-1)} - \sum_{l=1}^q \mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}}^{(l)} \right)^2. \quad (9)$$

We first discuss the estimation when the group memberships \mathcal{G} are given. In this case, the minimization of (9) is equivalent to solving $\frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \boldsymbol{\theta}_{g^{(l)}}^{(l)}} = \mathbf{0}$, $\frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \alpha_{g^{(1)} \dots g^{(q)}}} = 0$ for all $g^{(l)} \in [G_l]$

and $l \in [q]$. Recall that $N_{lg^{(l)}} = |\mathcal{R}_{g^{(l)}}^{(l)}|$ and let

$$\mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} = \left(\text{vec}\{(\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)})_{g^{(1)} \dots g^{(q)}}\}, \mathbf{1}_{N_{1g^{(1)}}} \otimes \cdots \otimes (\mathbf{X}_t^{(l)})^{\left(\mathcal{R}_{g^{(l)}}^{(l)}, \cdot\right)} \otimes \cdots \otimes \mathbf{1}_{N_{qg^{(q)}}} \right), \quad (10)$$

where $\mathbf{X}_t^{(l)} = (\mathbf{x}_{1t}^{(l)}, \dots, \mathbf{x}_{N_{lt}}^{(l)})^\top \in \mathbb{R}^{N_l \times p_l}$, and that

$$(\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)})_{g^{(1)} \dots g^{(q)}} = \left(\sum_{i_l=1}^{N_l} Y_{i_1, \dots, i_l, \dots, i_q, t-1} (a_{s, i_l}^{(l)} / n_{l s}) \right)_{i_1 \in \mathcal{R}_{g^{(1)}}^{(l)}, \dots, s \in \mathcal{R}_{g^{(l)}}^{(l)}, \dots, i_q \in \mathcal{R}_{g^{(q)}}^{(l)}}.$$

Then one can verify that

$$\frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \boldsymbol{\theta}_{g^{(l)}}^{(l)}} = \left(\sum_{t, g^{-(l)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \boldsymbol{\theta}_{g^{(l)}}^{(l)} \right. \\ \left. - \sum_{t, g^{-(l)}} \left\{ \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \left(\mathbb{Y}_{g^{(1)} \dots g^{(q)}, t} - \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)} \alpha_{g^{(1)} \dots g^{(q)}} - \sum_{m \neq l} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(m)} \boldsymbol{\theta}_{g^{(m)}}^{(m)} \right) \right\} \right), \quad (11)$$

where $\mathbb{Y}_{g^{(1)} \dots g^{(q)}, t} = \text{vec}\left\{ \mathcal{Y}_t^{\left(\mathcal{R}_{g^{(1)}}^{(1)}, \dots, \mathcal{R}_{g^{(q)}}^{(q)}\right)} \right\} \in \mathbb{R}^{N_{1g^{(1)}} \cdots N_{qg^{(q)}}}$, and the summation $\sum_{g^{-(l)}}$ is the simplified notation for $\sum_{g^{(1)}, \dots, g^{(l-1)}, g^{(l+1)}, \dots, g^{(q)}}$. Furthermore, it holds that

$$\frac{\partial Q(\boldsymbol{\xi}, \mathcal{G})}{\partial \alpha_{g^{(1)} \dots g^{(q)}}} = \sum_t \left\| \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)} \right\|^2 \alpha_{g^{(1)} \dots g^{(q)}} \\ - \sum_t \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^\top \left(\mathbb{Y}_{g^{(1)} \dots g^{(q)}, t} - \sum_{l=1}^q \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \boldsymbol{\theta}_{g^{(l)}}^{(l)} \right). \quad (12)$$

Equations (11) and (12) define a system of linear equations, whose solutions have the form $\widehat{\boldsymbol{\xi}} = \mathbf{M}^{-1} \mathbf{b}$ with $\mathbf{M} \in \mathbb{R}^{\{\sum_l G_l(p_l+1) + \prod_l G_l\} \times \{\sum_l G_l(p_l+1) + \prod_l G_l\}}$ and $\mathbf{b} \in \mathbb{R}^{\sum_l G_l(p_l+1) + \prod_l G_l}$,

where

$$\mathbf{M} = \begin{pmatrix} \mathbf{M}^{(1)} & \mathbf{M}^{(12)} & \mathbf{M}^{(13)} & \dots & \mathbf{M}^{(1q)} & \mathbf{M}^{(1\alpha)} \\ \mathbf{M}^{(21)} & \mathbf{M}^{(2)} & \mathbf{M}^{(23)} & \dots & \mathbf{M}^{(2q)} & \mathbf{M}^{(2\alpha)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{M}^{(q1)} & \mathbf{M}^{(q2)} & \dots & \dots & \mathbf{M}^{(q)} & \mathbf{M}^{(q\alpha)} \\ \mathbf{M}^{(1\alpha)\top} & \mathbf{M}^{(2\alpha)\top} & \dots & \dots & \mathbf{M}^{(q\alpha)} & \mathbf{M}^{(\alpha)} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}^{(1)} \\ \dots \\ \mathbf{b}^{(q)} \\ \mathbf{b}^\alpha \end{pmatrix}. \quad (13)$$

The specific terms in (13) are given in Appendix C. Subsequently, given the estimated parameters $\widehat{\boldsymbol{\xi}}$, we update the group memberships \mathcal{G}_l from $l = 1$ to q iteratively. Specifically, given $\boldsymbol{\xi}$ and $\mathcal{G}_{-l} \stackrel{\text{def}}{=} \{\mathcal{G} \setminus \mathcal{G}_l\}$, the \mathcal{G}_l is updated by

$$\widehat{g}_{i_l}^{(l)} \in \arg \min_{g_{i_l}^{(l)} \in [G_l]} \sum_{i_{-l}} \sum_{t=1}^T \left(Y_{i_1 i_2 \dots i_q, t} - \sum_{l=1}^q \lambda_{g_{i_l}^{(l)}}^{(l)} \sum_{k=1}^{N_l} \frac{a_{i_l k}^{(l)}}{n_{l i_l}} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)} \right. \\ \left. - \alpha_{g_{i_l}^{(l)} \dots g_{i_q}^{(q)}} Y_{i_1 i_2 \dots i_q, (t-1)} - \sum_{l=1}^q \mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}^{(l)}}^{(l)} \right)^2, \quad (14)$$

where $\sum_{i_{-l}} = \sum_{i_1=1}^{N_1} \dots \sum_{i_{l-1}=1}^{N_{l-1}} \sum_{i_{l+1}=1}^{N_{l+1}} \dots \sum_{i_q=1}^{N_q}$.

We summarize the algorithm in Algorithm A.2 in Appendix D, which consists of iterations of two major steps. The first step updates the estimated parameters given the group memberships, while the second step updates group memberships given the parameter estimates. We note that the update step in (14) requires input of memberships of other layers (i.e., \mathcal{G}_{-l}). Therefore, we adopt a sequential updating rule. During the k th iteration, for $l' < l$, we substitute $\mathcal{G}_{l'}$ by $\mathcal{G}_{l'}^{[k]}$, and for $l' > l$, we substitute $\mathcal{G}_{l'}$ by $\mathcal{G}_{l'}^{[k-1]}$ in the previous iteration. See (A.25) for detailed expression. Each step can be efficiently computed due to the simple analytical forms. This algorithm can be validated to converge to a local minimizer, of which the proof is given in Appendix E. We remark that whether it can converge to a global optimizer is a challenging question due to the nonconvexity of the objective function (9) (Murty and Kabadi, 1987; Lin et al., 2020). To increase the chance of finding the global optimizer, a set of initialization strategies is employed, which are described in detail in Algorithm A.3. We try multiple initial values and use the one with the minimum loss function as the best one, which guarantees a stable numerical performance. We leave the theoretical investigation for a global optimizer as an interesting future topic.

Remark 5. *In Algorithm A.2, STEP 2 achieves dimension reduction due to the calculation of group memberships. Note that in each layer, the update equation (A.25) only involves subject i_l , and does not depend on other inner-layer subjects. Thus, the memberships update equation (A.25) can be computed in parallel for each inner-layer subject (i_l in the l th layer), greatly increasing the computational speed. We show our computational cost under each sample size setting in Figure A.7. Even at the maximum network size configuration ($N_1 = 300, N_2 = 250$), total computational costs remain no more than 20 seconds across all time horizons. More discussion can be found in Appendix L.4.*

3.3 Selection of Group Numbers

To simplify the notations, denote $\underline{G} = (G_1, \dots, G_q)^\top$, and denote the corresponding true group numbers as \underline{G}_0 . We focus our later theoretical analysis on the case that G_l is finite

for $1 \leq l \leq q$. Write $\widehat{\boldsymbol{\xi}}(\underline{G}), \widehat{\boldsymbol{g}}(\underline{G})$ as the estimators when the group numbers are specified as \underline{G} . Then, we estimate \underline{G}_0 by utilizing the following information criterion:

$$\text{QIC}(\underline{G}) = \log\{Q(\widehat{\boldsymbol{\xi}}(\underline{G}), \widehat{\boldsymbol{g}}(\underline{G}))\} + \lambda(\underline{G}), \quad (15)$$

where $Q(\cdot, \cdot)$ is defined in (9), and $\lambda(\underline{G})$ is a penalty function. Then we estimate the group numbers by $\widehat{\underline{G}} \in \arg \min_{\underline{G}} \text{QIC}(\underline{G})$. In practice, we specify $\lambda(\underline{G}) = \kappa(\sum_l G_l)$, in which κ is a tuning parameter. Our theoretical analysis shows that if $T^{-1/2}(m + \sum_l \log N_l) \ll \kappa \ll c_{\text{gap}} c_{\pi}^q / (\prod_l G_l)$, the QIC can consistently select \underline{G}_0 . Here c_{gap} and c_{π} are group structures related values depending on the strength of model signals, which will be defined in our theoretical analysis in Section 4. In our numerical study, we specify $\kappa = \{C(\log T)T^{1/8}\}^{-1}$ with $C = 40$, which achieves reliable finite sample performances in all of our numerical studies.

4. Theoretical Properties

4.1 Estimation Consistency

Define $\boldsymbol{\Theta}_{i_1 \dots i_q} = (\boldsymbol{\theta}_{g_{i_1}^{(1)}}^{(1)\top}, \dots, \boldsymbol{\theta}_{g_{i_q}^{(q)}}^{(q)\top}, \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}})^{\top} \in \mathbb{R}^{\sum_l (p_l + 1) + 1}$ and $\boldsymbol{\Theta} = (\boldsymbol{\Theta}_{i_1 \dots i_q} : i_l \in [N_l], 1 \leq l \leq q)$ as a tensor of dimension $N_1 \times \dots \times N_q \times \{\sum_l (p_l + 1) + 1\}$. With $\boldsymbol{\Theta}$ we can rewrite the loss function (9) as follows

$$Q(\boldsymbol{\Theta}) = \sum_{i_1=1}^{N_1} \dots \sum_{i_q=1}^{N_q} \sum_{t=1}^T (Y_{i_1 \dots i_q, t} - \mathcal{X}_{i_1 \dots i_q, t}^{\top} \boldsymbol{\Theta}_{i_1 \dots i_q})^2 \stackrel{\text{def}}{=} \sum_{i_1=1}^{N_1} \dots \sum_{i_q=1}^{N_q} Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}), \quad (16)$$

where

$$\begin{aligned} \mathcal{X}_{i_1 \dots i_q, t} \stackrel{\text{def}}{=} & \left(\sum_{k=1}^{N_1} w_{i_1 k}^{(1)} Y_{ki_2 \dots i_q, (t-1)}, \mathbf{x}_{i_1 t}^{(1)\top}, \dots, \right. \\ & \left. \sum_{k=1}^{N_q} w_{i_q k}^{(q)} Y_{i_1 \dots i_{(q-1)} k, (t-1)}, \mathbf{x}_{i_q t}^{(q)\top}, Y_{i_1 \dots i_q, (t-1)} \right)^{\top} \in \mathbb{R}^{\sum_l (p_l + 1) + 1}. \end{aligned} \quad (17)$$

Denote by $\widehat{\boldsymbol{\Theta}} = (\widehat{\boldsymbol{\Theta}}_{i_1 \dots i_q} = (\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_1}^{(1)}}^{(1)\top}, \dots, \widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_q}^{(q)}}^{(q)\top}, \widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_q}^{(q)}})^{\top})$ the global minimizer of $Q(\boldsymbol{\Theta})$ with the estimated group memberships $\widehat{g}_i^{(l)}$ ($i \in [N_l], l \in [q]$), we define a pseudo distance as follows

$$d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}) = \sum_{l=1}^q \frac{1}{N_l} \sum_{i_l=1}^{N_l} \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)}}^{(l)}\|^2 + \frac{1}{\prod_l N_l} \sum_{i_1=1}^{N_1} \dots \sum_{i_q=1}^{N_q} |\widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}}|^2. \quad (18)$$

Intuitively, $d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta})$ measures the average distance between $\widehat{\boldsymbol{\Theta}}$ and $\boldsymbol{\Theta}$. Next, we establish the consistency of the global optimizer $\widehat{\boldsymbol{\Theta}}$ of the loss function (9) in this pseudo distance, for which we require the following definition and assumptions.

Definition 1. (*K-CONVEX CONCENTRATION*) Let $\mathbf{x} \in \mathbb{R}^n$ be a random vector. If for every 1-Lipschitz convex function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$, we have $E|\varphi(\mathbf{x})| < \infty$ and for every $t > 0$, it holds that $P\left(\left|\varphi(\mathbf{x}) - E\{\varphi(\mathbf{x})\}\right| \geq t\right) \leq 2\exp(-t^2/K^2)$, then \mathbf{x} is said to have the *K-convex concentration property*.

Assumption 1. (*PARAMETER SPACE*) The parameter satisfies that $\|\Theta\|_{\max} < \infty$.

Assumption 2. (*CONVEXITY*) Let $\Sigma_{i_1 \dots i_q} = E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top)$. We assume that $\tau_{\min} \stackrel{\text{def}}{=} \min_{i_1, \dots, i_q} \lambda_{\min}(\Sigma_{i_1 \dots i_q})$ is a positive constant.

Assumption 3. (*DISTRIBUTION OF NOISE TERM*) Assume $\varepsilon_{i_1 \dots i_q, t}$ is i.i.d across $i_l \in [N_l]$ for all $l \in [q]$ and $t \in [T]$. In addition, $\varepsilon_{i_1 \dots i_q, t}$ is a zero-mean sub-Gaussian variable with a scale factor $0 < \nu < \infty$, i.e., $E\{\exp(u\varepsilon_{i_1 \dots i_q, t})\} \leq \exp(\nu^2 u^2/2)$ for all $u \in \mathbb{R}$. Let $\varepsilon_{i_1 \dots i_q, t}$ be independent of $\{\mathcal{Y}_{s-1}, \mathbf{X}_s^{(1)}, \dots, \mathbf{X}_s^{(q)} : s \leq t\}$, where $\mathbf{X}_s^{(l)} = (\mathbf{x}_{i_l s}^{(l)} : i_l \in [N_l])^\top$.

Assumption 4. (*DISTRIBUTION OF COVARIATES*) Recall that $\mathbf{x}_{i_l t}^{(l)} \in \mathbb{R}^{p_l}$ is the covariate vector of the i_l th subject in the l th layer at time t . Assume $E(\mathbf{x}_{i_l t}^{(l)}) = \mathbf{0}$ for all $i_l \in [N_l]$ and $t \in [T]$. Let $\boldsymbol{\eta}_l \in \mathbb{R}^{p_l}$ be a constant vector satisfying $\|\boldsymbol{\eta}_l\| \leq c$ for $l \in [q]$, where c is a positive constant. Define $\mathbf{x}_t^{(l)\eta} = (\mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\eta}_l : i_l \in [N_l])^\top \in \mathbb{R}^{N_l}$. Assume $\mathbf{x}^{(l)\eta} = (\mathbf{x}_t^{(l)\eta\top} : 0 \leq t \leq T)^\top \in \mathbb{R}^{N_l(T+1)}$ satisfies the *K-convex concentration property* for some constant K according to Definition 1.

Assumption 5. (*STABILITY*) Denote $\mathbb{Y}_0 \stackrel{\text{def}}{=} \text{vec}(\mathcal{Y}_0) = \mathbf{0}$, and assume that

$$\max_{g^{(1)} \in [G_{1,0}], \dots, g^{(q)} \in [G_{q,0}]} \left| \sum_l \lambda_{g^{(l)}}^0 + \alpha_{g^{(1)} \dots g^{(q)}}^0 \right| \leq \kappa_{\max} < 1,$$

where $G_{l,0}$ is the true number of groups in the l th dimension, and κ_{\max} is a positive constant.

Assumption 6. (*GROUP DIFFERENCE*) Assume $\min_{g_1^{(l)} \neq g_2^{(l)}} \{\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_2^{(l)}}^{(l)0}\|^2 + \max_{g^{-(l)} \in [G_{-l}^0]} |\alpha_{g_1^{(l)} g^{-(l)}}^0 - \alpha_{g_2^{(l)} g^{-(l)}}^0|^2\} \geq c_{\text{gap}}$ holds for all $l \in [q]$, where $c_{\text{gap}} \gg T^{-1}(\sum_l \log N_l)^2$ as $T \rightarrow \infty$. Here $\{g^{-(l)} \in [G_{-l}^0]\}$ denotes $\{g^{(1)} \in [G_{1,0}], \dots, g^{(l-1)} \in [G_{l-1,0}], g^{(l+1)} \in [G_{l+1,0}], \dots, g^{(q)} \in [G_{q,0}]\}$.

Assumption 7. (*GROUP PROPORTION*) Let $\{g_i^{(l)0} : i_l \in [N_l]\}$ be non-random true membership sequences. Let $\pi_{g^{(l)}, N_l}^{(l)} = \sum_{i_l} I(g_i^{(l)0} = g^{(l)})/N_l$ for $g^{(l)} \in [G_{l,0}]$. Assume that $\min_{l \in [q], g^{(l)} \in [G_{l,0}]} \pi_{g^{(l)}, N_l}^{(l)} \geq c_\pi > 0$ for sufficiently large N_l , where c_π is a positive constant.

Assumption 1 requires the parameter space to be bounded. Assumption 2 ensures the convexity of the element-wise objective function, i.e., $Q_{i_1 \dots i_1}(\Theta_{i_1 \dots i_q})$, as a function of $\Theta_{i_1 \dots i_q}$ for sufficiently large T . This condition guarantees the unique solution of the local objective function for each node. and is crucial for establishing the consistency result for the metric (18). These conditions are also commonly used in the literature, see, Self and Liang (1987), Fan and Li (2001), Zou (2006), etc., for example.

Assumptions 3–4 concern about the distributions of the error term and covariates, respectively. Specifically, Assumption 3 requires the error term $\varepsilon_{i_1 \dots i_q, t}$ to be *i.i.d.* sub-Gaussian variables, which is widely used in high-dimensional time series literature (Wang et al., 2013; Lugosi and Mendelson, 2019; Fan et al., 2021). We also provide a weighted least squares estimation procedure with group-specific variances, i.e., $\text{var}(\varepsilon_{i_1 \dots i_q, t}) = \sigma_{g_{i_1}^{(1)0} \dots g_{i_q}^{(q)0}}$, in Section 5.2. Subsequently, Assumption 4 allows the covariates $\{\mathbf{x}_{i_t}^{(l)}\}$ to be correlated but satisfying the K -convex concentration property according to Definition 1. This assumption is employed to establish Hanson-Wright type inequality for dependent variables (Adamczak, 2015). Although this is a high-level condition, there are a variety of random variables satisfying Definition 1, as discussed in the following Remark 6. We further comment that the Assumptions 3–4 together imply that $\mathbf{v}^\eta \stackrel{\text{def}}{=}} (\mathbf{v}_t^\eta : 0 \leq t \leq T)^\top$ satisfies the K -convex concentration property for some constant K . Here $\mathbf{v}_t^\eta = (\mathbf{x}_t^{(1)\eta^\top}, \dots, \mathbf{x}_t^{(q)\eta^\top}, \mathbb{E}_t^\top)^\top \in \mathbb{R}^{\sum_l N_l + \prod_l N_l}$, where $\mathbb{E}_t = \text{vec}(\mathcal{E}_t)$.

Remark 6. *As discussed by Adamczak (2015), there are a variety of random vectors \mathbf{x} satisfying the K -convex concentration property in Definition 1. For example, (i) Any random vector \mathbf{x} with its elements x_i s independent for all i , and $|x_i| \leq 1$ a.s., satisfies Definition 1 (Talagrand, 1988); (ii) Any random vector \mathbf{x} with its elements in a bounded interval and geometrically strongly mixing satisfies Definition 1 (Samson, 2000). We refer to Adamczak (2015) for more detailed discussions.*

Next, Assumption 5 ensures the stability of the tensor-valued time series data as T goes to infinity, as defined in Lütkepohl (2005). Assumptions 6 and 7 are imposed on certain group properties. Assumption 6 assumes there is a gap between the true parameters of two different groups within any network. The condition is an extension of the same type of conditions assumed by the group panel data models with groups assigned on one dimension (Su et al., 2016; Ando and Bai, 2016; Zhang et al., 2019; Liu et al., 2020). In Assumption 6, special care is paid to the self-momentum parameter $\alpha_{g^{(1)} \dots g^{(q)}}$, indexed by the tensor dimension-wise group memberships, to ensure the parameter identifiability. Specifically, we require a min-max type condition for $\alpha_{g^{(1)} \dots g^{(q)}}$ in Assumption 6. Furthermore, instead of assuming $c_{\text{gap}} > c > 0$ by a positive constant c in existing literature (Bonhomme and Manresa, 2015; Su et al., 2016; Zhang et al., 2019; Liu et al., 2020), we allow $c_{\text{gap}} \rightarrow 0$ to study how this signal strength affects the theoretical properties. Lastly, Assumption 7 assumes a lower bound of group proportions. In the following, we establish the consistency of the pseudo distance.

Theorem 2. *Suppose $G_l \geq G_{l,0}$ for all $l \in [q]$, where $G_{l,0}$ is the true number of groups. In addition, assume Assumptions 1–5 hold. We have that $d(\widehat{\Theta}, \Theta^0) = O_p\{(\sum_l \log N_l)^2 T^{-1}\}$, where $\Theta^0 = (\Theta_{i_1 \dots i_q}^0 = (\boldsymbol{\theta}_{g_{i_1}^{(1)0}}^{(1)0^\top}, \dots, \boldsymbol{\theta}_{g_{i_q}^{(q)0}}^{(q)0^\top}, \alpha_{g_{i_1}^{(1)0} \dots g_{i_q}^{(q)0}}^0)^\top)$ is the true parameter for Θ .*

Theorem 2 implies that as long as we have $T \gg (\sum_l \log N_l)^2$, $\widehat{\Theta}$ is a consistent estimator for Θ^0 in the metric when G_l is possibly over-specified, i.e., $G_l \geq G_{l,0}$. This allows N_l to grow exponentially fast with \sqrt{T} , which is a mild condition compared to existing literatures (Su et al., 2016). As a result, it allows for moderately large N_l with respect to \sqrt{T} . Subsequently, we show in Theorem 3 that the QIC can consistently select the true group numbers.

Theorem 3. *Under Assumptions 1–7, and assume $\kappa = \lambda(\underline{G})/(\sum_l G_l)$ satisfies*

$$T^{-1}(\sum_l \log N_l)^2 \ll \kappa \ll c_{\text{gap}}/(\prod_l G_l). \quad (19)$$

Then we have $P(\widehat{G}_1 = G_{1,0}, \dots, \widehat{G}_q = G_{q,0}) \rightarrow 1$ as $\min\{N_1, \dots, N_q, T\} \rightarrow \infty$ and $T^{-1}(\sum_l \log N_l)^2 \rightarrow 0$.

Theorem 3 implies that if we set κ to satisfy (19), then we can consistently estimate the true group numbers. Denote $G_{-l} = (G_k : k \neq l)^\top \in \mathbb{R}^{q-1}$, we need $\kappa \gg T^{-1}(\sum_l \log N_l)^2$ to ensure that $\text{QIC}(\underline{G}) > \text{QIC}(\underline{G}_0)$ for the over-fitting case, i.e., $G_l > G_{l,0}$ and $G_{-l} \geq G_{-l,0}$ for any $l \in [q]$. Here $G_{-l} \geq G_{-l,0}$ means $G_k \geq G_{k,0}$ for any $k \neq l$. Conversely, we need $\kappa \ll c_{\text{gap}}/(\prod_l G_l)$ to guarantee that $\text{QIC}(\underline{G}) > \text{QIC}(\underline{G}_0)$ for the under-fitting case, i.e., there exists an l with $G_l < G_{l,0}$. When both conditions are met, we can obtain $\widehat{G}_l = G_{l,0}$ for all $l \in [q]$ with a probability approaching 1. In the next subsection, we further discuss the results of node-wise parameter estimation and the strong group membership consistency.

4.2 Membership Estimation Consistency and Asymptotic Normality

As we stated before, the metric in (18) measures the average estimation error between $\widehat{\Theta}$ and Θ^0 . Therefore, the result in Theorem 2 is not sufficient to imply the parameter consistency for each node. To this end, we derive the following node-wise parameter consistency result, which will be crucial to building the strong membership estimation consistency later.

Proposition 4. *Under Assumptions 1–5, when $G_l \geq G_{l,0}$ for all $l \in [q]$, we have*

$$\begin{aligned} & \sup_{i_l} \left\{ \left\| \widehat{\theta}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \theta_{g_{i_l}^{(l)0}}^{(l)0} \right\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m=1}^{N_m} \left| \widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)0} \dots g_{i_q}^{(q)0}} \right|^2 \right\} \\ & = O_p \left\{ T^{-1}(\sum_l \log N_l)^2 \right\}. \end{aligned} \quad (20)$$

Recall that $\mathcal{G}_l = (g_{i_l}^{(l)} : 1 \leq i_l \leq N_l)^\top$, and let \mathcal{G}_l^0 denote the corresponding true memberships. Equation (20) establishes uniform node-wise parameter estimation consistency, which is crucial for achieving the following strong consistency of membership estimation for $\widehat{\mathcal{G}}_l$ when $G_l \geq G_{l,0}$ for all $l \in [q]$. Denote the node set $\widehat{\mathcal{R}}_{\widetilde{g}^{(l)}}^{(l)} = \{i_l : \widehat{g}_{i_l}^{(l)} = \widetilde{g}^{(l)}\}$ for $\widetilde{g}^{(l)} \in [G_l]$ and $\mathcal{R}_{g^{(l)}}^{(l)0} = \{i_l : g_{i_l}^{(l)0} = g^{(l)}\}$ for $g^{(l)} \in [G_{l,0}]$, where $\widehat{g}_{i_l}^{(l)}$ is obtained by (14) when $\widehat{\xi}$ is specified.

Theorem 5. *(Strong Consistency of Membership Estimation) Under Assumptions 1–7, and suppose $G_l \geq G_{l,0}$ for all $l \in [q]$. Then for each estimated group $\widetilde{g}^{(l)} \in [G_l]$, there exists a true group $g^{(l)} \in [G_{l,0}]$, such that $\widehat{\mathcal{R}}_{\widetilde{g}^{(l)}}^{(l)} \subset \mathcal{R}_{g^{(l)}}^{(l)0}$ with probability tending to 1.*

As implied by Theorem 5, the true groups are split into subgroups instead of joining into new groups when $G_l \geq G_{l,0}$ for all $l \in [q]$, as shown in Figure A.6. Define $\widehat{\mathcal{G}}_l = (\widehat{g}_{i_l}^{(l)} : i_l \in [N_l])^\top \in \mathbb{R}^{N_l}$ as the estimated membership vectors for all $l \in [q]$. Particularly, when $G_l = G_{l,0}$, we can show that $\widehat{\mathcal{G}}_l = \mathcal{G}_l^0$ holds with probability tending to

1 under certain label permutations. We would like to remark that establishing the group consistency result is non-trivial due to several challenges under our framework. As shown in the GTNAR model (3), group memberships from different dimensions are involved and tangled in an additive form. In addition, the network dependence structure creates extra difficulty in establishing the node-wise membership consistency result. To establish the membership consistency, for instance, in Zhu et al. (2023), a refinement procedure is further required for obtaining this property. In Su et al. (2023), a two step estimation procedure is used to achieve group membership estimation consistency when the true group number is known. Specifically, their two-step procedure establishes node-wise parameter convergence and consistent group membership identification sequentially, whereas we allow simultaneous estimation of group membership and corresponding group-wise parameters within a unified framework. In contrast to their work, a further investigation of our additive model form enables us to address these challenges and establish the group membership consistency result for stochastically estimated $\widehat{\mathcal{G}}_l$ directly, which is critically important for later statistical inference. Furthermore, let $\widehat{\boldsymbol{\xi}}^{\text{or}}$ be the oracle estimator when the true group memberships \mathcal{G}_l^0 s are known for all $l \in [q]$. Then the oracle property holds that $\widehat{\boldsymbol{\xi}} = \widehat{\boldsymbol{\xi}}^{\text{or}}$ with probability tending to 1, where $\widehat{\boldsymbol{\xi}}$ is the estimator for $\boldsymbol{\xi}^0 = (\boldsymbol{\theta}^{(1)0\top}, \dots, \boldsymbol{\theta}^{(q)0\top}, \text{vec}(\boldsymbol{\alpha}^0)^\top)^\top$. The results are presented in the following Corollary.

Corollary 6. *Under Assumptions 1–7, and assume $G_l = G_{l,0}$ for all $l \in [q]$. Then under label permutations, we have*

$$\lim_{\min\{N_1, \dots, N_q, T\} \rightarrow \infty} P\left(\widehat{\mathcal{G}}_1 = \mathcal{G}_1^0, \dots, \widehat{\mathcal{G}}_q = \mathcal{G}_q^0\right) \rightarrow 1, \quad (21)$$

$$\lim_{\min\{N_1, \dots, N_q, T\} \rightarrow \infty} P\left(\widehat{\boldsymbol{\xi}} = \widehat{\boldsymbol{\xi}}^{\text{or}}\right) \rightarrow 1. \quad (22)$$

The results in Corollary 6 imply that $\widehat{\boldsymbol{\xi}}$ is asymptotically equivalent to $\widehat{\boldsymbol{\xi}}^{\text{or}}$. Therefore, to derive the asymptotic distribution of $\widehat{\boldsymbol{\xi}}$, it is sufficient to investigate the case for $\widehat{\boldsymbol{\xi}}^{\text{or}}$. We establish the asymptotic normality result in the following theorem.

Theorem 7. *Assume Assumptions 1–7, $G_l = G_{l,0}$ and that there exists n , such that $c_1 n \leq \min_l N_l \leq \max_l N_l \leq c_2 n$ for some constants $c_1, c_2 > 0$. Define $\mathbf{M}_{nT}^0 = n^{-q} T^{-1} E(\mathbf{M})$ and assume $\mathbf{M}^0 = \lim_{\min\{n, T\} \rightarrow \infty} \mathbf{M}_{nT}^0$ exists, where \mathbf{M} is given in (13). Assume $\lambda_{\min}(\mathbf{M}^0) \geq \tau > 0$ for a positive constant τ . Let $k = \sum_l \{G_l(p_l + 1)\} + \prod_l G_l$. Then for any $\boldsymbol{\eta} \in \mathbb{R}^k$ with $\|\boldsymbol{\eta}\| = 1$ we have*

$$n^{q/2} T^{1/2} \boldsymbol{\eta}^\top (\widehat{\boldsymbol{\xi}} - \boldsymbol{\xi}^0) \rightarrow_d N\left(\mathbf{0}, \sigma^2 \boldsymbol{\eta}^\top (\mathbf{M}^0)^{-1} \boldsymbol{\eta}\right). \quad (23)$$

Theorem 7 establishes the asymptotic normality of the estimator. Specifically, the convergence rates of $\widehat{\boldsymbol{\theta}}^{(l)}$ is $\sqrt{n^q T}$ for all $l \in [q]$. Using (23), we can conduct the statistical inference.

5. Model Extensions

The GTNAR model can be flexibly extended to accommodate an interactive model setting and group-specific error variances. In this section, we present the corresponding estimation and establish its statistical guarantees.

5.1 Mixed GTNAR Model

GTNAR model can be extended to incorporating the interactive network effects as

$$\begin{aligned} \mathcal{Y}_t &= (\mathcal{Y}_{t-1} \times_{l=1}^q \mathbf{W}^{(l)}) \times_{l=1}^q \mathbf{\Gamma}^{(l)} \\ &+ \sum_l (\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)}) \times_l \mathbf{L}^{(l)} + \mathbf{A} \odot \mathcal{Y}_{t-1} + \sum_{l=1}^q \boldsymbol{\beta}_{X_{l,t}}^{(l)} \circ_{k \neq l} \mathbf{1}_{N_k} + \mathbf{E}_t, \end{aligned} \quad (24)$$

which we refer to as the Mixed GTNAR model. The first term represents the interactive network effects with $\mathbf{\Gamma}^{(l)} = \text{diag}(\gamma_{g_i}^{(l)} : i_l \in [N_l]) = \text{diag}(\gamma^{(l)})$. Specifically, the interactive network coefficients are involved in (24) in a multiplicative form across different layers. Take the case of $q = 2$ as an example, model (24) can be expressed as

$$\begin{aligned} \mathbf{Y}_t &= (\mathbf{\Gamma}^{(1)} \mathbf{W}^{(1)}) \mathbf{Y}_{t-1} (\mathbf{W}^{(2)} \mathbf{\Gamma}^{(2)}) + (\mathbf{L}^{(1)} \mathbf{W}^{(1)}) \mathbf{Y}_{t-1} + \mathbf{Y}_{t-1} (\mathbf{W}^{(2)} \mathbf{L}^{(2)}) \\ &+ \mathbf{A} \odot \mathbf{Y}_{t-1} + \boldsymbol{\beta}_{X_{1,t}}^{(1)} \mathbf{1}_{N_2}^\top + \mathbf{1}_{N_1} \boldsymbol{\beta}_{X_{2,t}}^{(2)\top} + \mathbf{E}_t. \end{aligned} \quad (25)$$

Note that the multiplication form in $\mathbf{\Gamma}^{(1)} \mathbf{W}^{(1)} \mathbf{Y}_{t-1} \mathbf{W}^{(2)} \mathbf{\Gamma}^{(2)}$ causes the parameter identification issue. Following the convention (Chen et al., 2021), we set $\|\mathbf{\Gamma}^{(1)}\|_F = 1$ to guarantee the interactive network effects $\mathbf{\Gamma}^{(1)}$ and $\mathbf{\Gamma}^{(2)}$ to be identifiable with sign flips. To estimate model (25), we apply an iterative least squares method, and the detailed estimation algorithms are given in Appendix G.1. In the following, we present the theoretical properties for the estimator of model (24). Note that one can re-write the first interactive term in (25) as $(\boldsymbol{\gamma}^{(1)} \boldsymbol{\gamma}^{(2)\top}) \odot (\mathbf{W}^{(1)} \mathbf{Y}_{t-1} \mathbf{W}^{(2)})$, hence $\boldsymbol{\gamma}^{(1)} \boldsymbol{\gamma}^{(2)\top}$ plays a similar role as \mathbf{A} . As a result, we can borrow the proof idea for dealing with the autoregressive matrix \mathbf{A} in our original model (4) with $q = 2$ to establish the consistency result for the Mixed GTNAR model. We first introduce some necessary conditions. For notational simplicity, denote $\boldsymbol{\Theta}_{ij} = (\boldsymbol{\theta}_{g_i^{(1)}}^{(1)\top}, \boldsymbol{\theta}_{g_j^{(2)}}^{(2)\top}, \alpha_{g_i^{(1)} g_j^{(2)}}, \gamma_{g_i^{(1)} g_j^{(2)}})^\top \in \mathbb{R}^{p_1+p_2+4}$, and denote $\boldsymbol{\Theta} = (\boldsymbol{\Theta}_{ij} : i \in [N_1], j \in [N_2]) \in \mathbb{R}^{N_1 \times N_2 \times (p_1+p_2+4)}$. Let $\boldsymbol{\Theta}^0 = (\boldsymbol{\Theta}_{ij}^0 = (\boldsymbol{\theta}_{g_i^{(1)}}^{(1)0\top}, \boldsymbol{\theta}_{g_j^{(2)}}^{(2)0\top}, \alpha_{g_i^{(1)} g_j^{(2)}}^0, \gamma_{g_i^{(1)} g_j^{(2)}}^0)^\top)$ be the true parameter for $\boldsymbol{\Theta}$. Define the pseudo distance as

$$\begin{aligned} d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}) &= \frac{1}{N_1} \sum_{i=1}^{N_1} \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_i^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_i^{(1)}}^{(1)}\|^2 + \frac{1}{N_2} \sum_{j=1}^{N_2} \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_j^{(2)}}^{(2)} - \boldsymbol{\theta}_{g_j^{(2)}}^{(2)}\|^2 \\ &+ \frac{1}{N_1 N_2} \sum_{i=1}^{N_1} \sum_{j=2}^{N_2} \left\{ |\widehat{\alpha}_{\widehat{g}_i^{(1)} \widehat{g}_j^{(2)}} - \alpha_{g_i^{(1)} g_j^{(2)}}|^2 + |\widehat{\gamma}_{\widehat{g}_i^{(1)} \widehat{g}_j^{(2)}}^{(1)} \widehat{\gamma}_{\widehat{g}_j^{(2)}}^{(2)} - \gamma_{g_i^{(1)} g_j^{(2)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)}|^2 \right\}. \end{aligned} \quad (26)$$

In the above distance, we note that the interactive effect $\widehat{\gamma}_{\widehat{g}_i^{(1)} \widehat{g}_j^{(2)}}^{(1)} \widehat{\gamma}_{\widehat{g}_j^{(2)}}^{(2)}$ plays a similar role as $\widehat{\alpha}_{\widehat{g}_i^{(1)} \widehat{g}_j^{(2)}}$ in $d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta})$. Hence, we can borrow the proof idea of Theorem 5 to establish the estimation consistency result with $d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta})$ for model (25). The required assumptions are given in Appendix G.2. We state the formal theoretical results in the following.

Theorem 8. *Under Assumptions 2–4 and A.11–A.13, the following conclusions hold.*

- (i) *Suppose $G_l \geq G_{l,0}$ for $l = 1, 2$, then $d(\widehat{\Theta}, \Theta^0) = O_p\{T^{-1}(\log(N_1 N_2))^2\}$.*
- (ii) *Assume $\kappa_1 = \lambda(\underline{G})/(G_1 + G_2)$ satisfies*

$$T^{-1}(\log(N_1 N_2))^2 \ll \kappa_1 \ll c_{\text{gap}}/(G_1 G_2). \quad (27)$$

Then $P(\widehat{G}_1 = G_{1,0}, \widehat{G}_2 = G_{2,0}) \rightarrow 1$ as $\{N_1, N_2, T\} \rightarrow \infty$ and $T^{-1}(\log(N_1 N_2))^2 \rightarrow 0$.

- (iii) *Further suppose Assumption 7 and A.14 hold, and $G_l \geq G_{l,0}$ for $l = 1, 2$. Then for any estimated group $\widehat{g}^{(1)} \in [G_1]$ and $\widehat{g}^{(2)} \in [G_2]$, there exists true groups $g^{(1)} \in [G_{1,0}]$ and $g^{(2)} \in [G_{2,0}]$, such that*

$$P\left\{\widehat{\mathcal{R}}_{\widehat{g}^{(1)}}^{(1)} \in \mathcal{R}_{g^{(1)0}}^{(1)0}\right\} \rightarrow 1, \quad P\left\{\widehat{\mathcal{R}}_{\widehat{g}^{(2)}}^{(2)} \in \mathcal{R}_{g^{(2)0}}^{(2)0}\right\} \rightarrow 1,$$

where $\widehat{\mathcal{R}}_{\widehat{g}^{(l)}}^{(l)} = \{i_l : \widehat{g}_i^{(l)} = \widehat{g}^{(l)}\}$ for $l = 1, 2$.

In Theorem 8 (i), we first establish the estimation consistency based on the pseudo distance $d(\widehat{\Theta}, \Theta)$ when the group numbers are possibly over-specified. It shows that when the interactive model is correctly specified or over-specified, the parameter $\widehat{\Theta}$ is a consistent estimator for Θ^0 as long as $T \gg (\log(N_1 N_2))^2$. Next in (ii), we show that the QIC selection method in (A.30) leads to consistent group number estimators. When the tuning parameter κ_1 satisfies the condition (27), we can consistently estimate the group numbers. The lower bound and upper bound for κ_1 are applied for the results in over- and under-specified situations, respectively. Subsequently, we establish the strong group membership estimation consistency in (iii). The proof of Theorem 8 is given in Appendix G.4. Since the interactive parameter $\gamma_{g_i^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)}$ plays a similar role as $\alpha_{g_i^{(1)} g_j^{(2)}}$ in model (25), the proofs exhibit analogous roadmap to the proofs of Theorem 2–5 in Section 4. However, the new interactive term in $\mathcal{X}_{ij,t}$ defined in (A.31) requires additional technical Lemmas 13–16 and Lemma 18, which are all crucial to the proof of the theorem.

To demonstrate the finite sample performance, we conduct several simulation experiments, which can be found in Appendix G.3.

5.2 Weighted Least Squares Estimation with Group-specific Error Variances

Next, we design a weighted least squares estimator for the case where the error terms have group-specific variances when $q = 2$. Several simulation studies are presented in Appendix B.4.5.

5.2.1 ESTIMATION PROCEDURE

In the case where the error term has group-specific variance, the GTNAR model can be modified as

$$\begin{aligned} Y_{ij,t} = & \lambda_{g_i^{(1)}}^{(1)} \sum_{k=1}^{N_1} w_{1ik} Y_{kj,(t-1)} + \lambda_{g_j^{(2)}}^{(2)} \sum_{k=1}^{N_2} Y_{ik,(t-1)} w_{2kj} + \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)} \\ & + \mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} + \varepsilon_{ij,t}, \end{aligned} \quad (28)$$

where $\varepsilon_{ij,t}$ satisfies that $E(\varepsilon_{ij,t}) = 0$ and $\text{var}(\varepsilon_{ij,t}) = \sigma_{g_i^{(1)0}g_j^{(2)0}}^2$, and $g_i^{(1)0}$ are $g_j^{(2)0}$ are true memberships. Under this situation, a weighted least squares (WLS) estimation procedure can be introduced (Carroll and Cline, 1988; Shao, 1989). Specifically, similar to the calculation for (A.7), a weighted version for the corresponding terms can be written as

$$\tilde{\mathbf{M}} = \begin{pmatrix} \tilde{\mathbf{M}}^{(1)} & \tilde{\mathbf{M}}^{(12)} & \tilde{\mathbf{M}}^{(1\alpha)} \\ \tilde{\mathbf{M}}^{(12)\top} & \tilde{\mathbf{M}}^{(2)} & \tilde{\mathbf{M}}^{(2\alpha)} \\ \tilde{\mathbf{M}}^{(1\alpha)\top} & \tilde{\mathbf{M}}^{(2\alpha)\top} & \tilde{\mathbf{M}}^\alpha \end{pmatrix}, \quad \tilde{\mathbf{b}} = \begin{pmatrix} \tilde{\mathbf{b}}^{(1)} \\ \tilde{\mathbf{b}}^{(2)} \\ \tilde{\mathbf{b}}^\alpha \end{pmatrix}. \quad (29)$$

Take the first term $\tilde{\mathbf{M}}^{(1)}$ as an example, we set

$$\tilde{\mathbf{M}}^{(1)} = \text{diag}\{\tilde{\mathbf{M}}_{g^{(1)}}^{(1)} : g^{(1)} \in [G_1]\} \in \mathbb{R}^{G_1(p_1+1) \times G_1(p_1+1)}, \quad (30)$$

$$\tilde{\mathbf{M}}_{g^{(1)}}^{(1)} = \sum_{t, g^{(2)}} \sigma_{g^{(1)g^{(2)}}}^{-2} \mathbb{X}_{g^{(1)g^{(2)}, t}}^\top \mathbb{X}_{g^{(1)g^{(2)}, t}} \quad (31)$$

with additional weights $\sigma_{g^{(1)g^{(2)}}}^{-2}$. Other terms in (29) can be calculated similarly. Then, the WLS estimator is obtained as $\tilde{\boldsymbol{\theta}} = \tilde{\mathbf{M}}^{-1}\tilde{\mathbf{b}}$.

However, there are still unknown parameters (e.g., $\sigma_{g^{(1)g^{(2)}}}^2$) in $\tilde{\boldsymbol{\theta}}$. Therefore, to obtain a feasible WLS estimator, we first estimate $\sigma_{g^{(1)g^{(2)}}}^2$ as follows

$$\hat{\sigma}_{g^{(1)g^{(2)}}}^2 = \frac{1}{N_{1g^{(1)}}N_{2g^{(2)}}T} \sum_{i \in \hat{\mathcal{R}}_{g^{(1)}}^{(1)}} \sum_{j \in \hat{\mathcal{R}}_{g^{(2)}}^{(2)}} \sum_t (Y_{ij,t} - \mathcal{X}_{ij,t}^\top \hat{\boldsymbol{\Theta}}_{ij})^2 \quad (32)$$

and $\hat{\boldsymbol{\Theta}}_{ij}$ is the estimator obtained by Algorithm A.2. This allows us to obtain the WLS estimator as

$$\tilde{\boldsymbol{\theta}}^w = \tilde{\mathbf{M}}^{-1}\tilde{\boldsymbol{\delta}}, \quad (33)$$

where $\tilde{\mathbf{M}}$ and $\tilde{\boldsymbol{\delta}}$ are obtained by substituting the estimators $\hat{\sigma}_{g^{(1)g^{(2)}}}^2$ and $\hat{\boldsymbol{\Theta}}$ in $\tilde{\mathbf{M}}$ and $\tilde{\mathbf{b}}$ of (29). In the next subsection, we establish the theoretical properties for $\tilde{\boldsymbol{\theta}}^w$.

5.2.2 THEORETICAL PROPERTIES

We first show that when the group-specific error variances exist, the membership estimation consistency in Theorem 5 still holds. Technically, we need the Assumption A.10 in Appendix B.4.1 instead of the Assumption 3 in Section 4. We state the formal theoretical results in Theorem 9.

Theorem 9. *Assume the group-specific error variances exist and the model is formed as (28). Suppose that Assumptions 1, 2, 4-7, A.10 hold. When $G_1 = G_{1,0}, G_2 = G_{2,0}$, we have,*

(i) *under label permutation,*

$$\lim_{\min(N_1, N_2, T) \rightarrow \infty} P\left(\hat{\mathcal{G}}^{(1)} = \mathcal{G}^{(1)0}, \hat{\mathcal{G}}^{(2)} = \mathcal{G}^{(2)0}\right) \rightarrow 1, \quad (34)$$

$$\lim_{\min(N_1, N_2, T) \rightarrow \infty} P(\tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}^{or}) \rightarrow 1. \quad (35)$$

Further assume that there exists n , such that $c_1 n \leq \min_l N_l \leq \max_l N_l \leq c_2 n$ and suppose that $\sqrt{T} \gg \log(G_1 G_2 n^2)$,

(ii) it holds that

$$\sup_{g^{(1)}, g^{(2)}} \left| \frac{1}{\hat{\sigma}_{g^{(1)}g^{(2)}}^2} - \frac{1}{\sigma_{g^{(1)}g^{(2)}}^2} \right| = o_p(1).$$

(iii) Assume $\lambda_{\min}(\mathbf{M}^0) \geq \tau > 0$ for a positive constant τ . Let $s = G_1(p_1 + 1) + G_2(p_2 + 1) + G_1 G_2$, then for any $\boldsymbol{\eta} \in \mathbb{R}^s$ with $\|\boldsymbol{\eta}\| = 1$, we have

$$(n\sqrt{T})^{-1} \boldsymbol{\eta}^\top (\tilde{\boldsymbol{\theta}}^w - \boldsymbol{\theta}^0) \rightarrow_d N(\mathbf{0}, \boldsymbol{\eta}^\top (\tilde{\mathbf{M}}^0)^{-1} \boldsymbol{\eta}), \quad (36)$$

where $\tilde{\mathbf{M}}^0 = \lim_{n, T \rightarrow \infty} (n^2 T)^{-1} E(\tilde{\mathbf{M}})$, and $\tilde{\mathbf{M}}$ is defined in (29).

The proof of Theorem 9 is provided in Appendix B.4.2–B.4.4. Conclusion (i) shows that in spite of the group-specific error variance, the group membership maintains strong consistency under nearly the same conditions. By conclusion (i), we have $\hat{\mathcal{G}}^{(1)} = \mathcal{G}^{(1)0}$ and $\hat{\mathcal{G}}^{(2)} = \mathcal{G}^{(2)0}$ with probability tending to 1, hence we next treat the group memberships as known, denoted as $\mathcal{R}_{g^{(1)}}^{(1)}$ and $\mathcal{R}_{g^{(2)}}^{(2)}$ in the technical proof. Note that we substitute $\sigma_{g^{(1)}g^{(2)}}^{-2}$ in (31) to be its estimator $\hat{\sigma}_{g^{(1)}g^{(2)}}^{-2}$. Hence, before we show the theoretical convergence properties of $\tilde{\boldsymbol{\theta}}^w$, we present the estimation consistency for $\hat{\sigma}_{g^{(1)}g^{(2)}}^{-2}$ in (ii), which is a critical result for later analysis. Conclusion (ii) shows that when $\sqrt{T} \gg \log(G_1 G_2 n^2)$, the sample weight $\hat{\sigma}_{g^{(1)}g^{(2)}}^{-2}$ is a consistent estimator of $\sigma_{g^{(1)}g^{(2)}}^{-2}$. Next, we establish the asymptotic normality for the WLS estimator $\tilde{\boldsymbol{\theta}}^w = \tilde{\mathbf{M}}^{-1} \tilde{\boldsymbol{\delta}}$ in (iii). In practice, one can calculate the estimator for $\tilde{\mathbf{M}}^0$ as $\tilde{\mathbf{M}}_{nT} = (n^2 T)^{-1} \tilde{\mathbf{M}}$, where $\tilde{\mathbf{M}}$ is obtained by plugging in the error variance estimator $\hat{\sigma}_{g^{(1)}g^{(2)}}^{-2}$ calculated by (32) in $\tilde{\mathbf{M}}$ defined in (29). In addition, we evaluate the finite sample performance of the WLS estimator $\tilde{\boldsymbol{\theta}}^w$ in Appendix B.4.5.

6. Simulation Study

6.1 Model Settings

To evaluate the finite sample performance of the GTNAR method, we conduct several simulation studies in this section for $q = 2$ (i.e., two networks). Three different scenarios for $G_{l,0}$ are considered. In each scenario, the node memberships $g_i^{(l)}$ s are sampled from the multinomial distribution with probability $\boldsymbol{\pi}_l = \{\pi_{g^{(l)}}^{(l)} = G_{l,0}^{-1} : g^{(l)} = 1, \dots, G_{l,0}\}$, $l \in [q]$. The dimension of exogenous covariates is set as $p_l = 3$, and the corresponding true parameters are shown in Table 1. For all scenarios, the covariates $\mathbf{x}_{it}^{(l)}$, $l \in [q]$ are generated from multivariate normal distribution $N(\mathbf{0}, \mathbf{I}_{p_l})$. Subsequently, we generate the noise term $\varepsilon_{i_1 i_2, t}$ from $N(0, 1)$ independently. For each scenario, the following two network structures are considered.

Table 1: True parameters for different scenarios.

Scenario 1		Scenario 2		Scenario 3	
$G_{1,0}$	$G_{2,0} = 2$	$G_{1,0} = 3$	$G_{2,0} = 2$	$G_{1,0} = 3$	$G_{2,0} = 3$
$\lambda_{g^{(1)}}^{(1)}$	$\lambda_{g^{(2)}}^{(2)}$	$\lambda_{g^{(1)}}^{(1)}$	$\lambda_{g^{(2)}}^{(2)}$	$\lambda_{g^{(1)}}^{(1)}$	$\lambda_{g^{(2)}}^{(2)}$
(0.15, 0.2)	(0.25, 0.4)	(0.15, 0.2, 0.3)	(0.25, 0.3)	(0.15, 0.2, 0.3)	(0.25, 0.3, 0.4)
$\zeta_{g^{(1)}}^{(1)}$	$\delta_{g^{(2)}}^{(2)}$	$\zeta_{g^{(1)}}^{(1)}$	$\delta_{g^{(2)}}^{(2)}$	$\zeta_{g^{(1)}}^{(1)}$	$\delta_{g^{(2)}}^{(2)}$
$\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.3 & 0.35 \\ 0.2 & -0.25 & 0.32 \end{pmatrix}$	$\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \\ 0.24 & 0.3 & -0.32 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.3 & 0.35 \\ 0.2 & -0.25 & 0.32 \end{pmatrix}$	$\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \\ 0.24 & 0.30 & -0.32 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.3 & 0.35 \\ 0.2 & -0.25 & 0.32 \\ 0.1 & -0.2 & 0.2 \end{pmatrix}$
$\alpha_{g^{(1)g^{(2)}}}$		$\alpha_{g^{(1)g^{(2)}}}$		$\alpha_{g^{(1)g^{(2)}}}$	
$\begin{pmatrix} -0.2 & 0.3 \\ -0.18 & 0.35 \end{pmatrix}$		$\begin{pmatrix} -0.2 & 0.3 \\ -0.18 & 0.35 \\ -0.15 & 0.28 \end{pmatrix}$		$\begin{pmatrix} -0.2 & 0.3 & 0.4 \\ -0.18 & 0.35 & 0.4 \\ -0.15 & 0.28 & 0.2 \end{pmatrix}$	

Example 1. (Stochastic Block Model, SBM) The first type of network is the stochastic block model, in which nodes in the same block (group) are assigned with a higher probability to be connected. In the l th network, based on the group memberships $g_i^{(l)}$ s and following the setting of Nowicki and Snijders (2001), we set $P(a_{ij}^{(l)} = 1) = 20/N_l$ when the i th and the j th node are in the same group, and otherwise we set $P(a_{ij}^{(l)} = 1) = 2/N_l$.

Example 2. (Power-Law Distribution Network) The second type of network is generated from a power-law distribution following Clauset et al. (2009). For the i th node in the l th network, its in-degree $d_i^{(l)} = \sum_{j=1}^N a_{ji}^{(l)}$ is assumed to be power-law distributed. Specifically, we first generate $\tilde{d}_i^{(l)}$ with a probability $P(\tilde{d}_i^{(l)} = k) \propto k^{-2.5}$, and set $d_i^{(l)} = 4\tilde{d}_i^{(l)}$. Then, $d_i^{(l)}$ followers of the i th node are randomly selected to construct the adjacency matrix. As a result, the adjacency matrix is not symmetric, which implies a directed network.

6.2 Performance Measure and Simulation Results

We first introduce the model performance measure and then present the simulation results. We set the network sizes $(N_1, N_2) \in \{(100, 80), (200, 150), (300, 250)\}$. The time length is set to be $T \in \{20, 40\}$. For each scenario, we repeat the experiments for $R = 500$ times. The networks are fixed throughout all replicates under one setting. In the initialization, we use 3 trials for each clustering type, as described in Algorithm A.3 of Appendix F. Denote the estimated parameters in the r th replicate as $\hat{\lambda}_{g^{(1)}}^{(1)[r]}$, $\hat{\lambda}_{g^{(2)}}^{(2)[r]}$, $\hat{\zeta}_{g^{(1)}}^{(1)[r]}$, $\hat{\zeta}_{g^{(2)}}^{(2)[r]}$, $\hat{\alpha}_{g^{(1)g^{(2)}}}$ and the estimated group number as $\hat{G}_1^{[r]}$ and $\hat{G}_2^{[r]}$.

6.2.1 ESTIMATION WHEN $G_1 = G_{1,0}$ AND $G_2 = G_{2,0}$

We first evaluate the estimation accuracy when the group numbers are correctly specified. Take $\boldsymbol{\lambda}^{(1)} = (\lambda_1^{(1)}, \dots, \lambda_{G_1}^{(1)})^\top$ for example. Denote $\hat{\boldsymbol{\lambda}}^{(1)[r]}$ as the estimator of $\boldsymbol{\lambda}^{(1)0} = (\lambda_1^{(1)0}, \dots, \lambda_{G_1}^{(1)0})^\top$ in the r th replicate. To evaluate the estimation accuracy, we calculate the root mean squared error (RMSE) as $\text{RMSE}_{\boldsymbol{\lambda}^{(1)}} = \{R^{-1} \sum_{r=1}^R (\|\hat{\boldsymbol{\lambda}}^{(1)[r]} - \boldsymbol{\lambda}^{(1)0}\|^2)\}^{1/2}$. Next, to gauge the performance of the statistical inference, we construct the 95% confidence interval for each parameter. For example, denote the estimated standard error of $\lambda_{g^{(1)}}^{(1)}$ as $\widehat{\text{SE}}_{\lambda_{g^{(1)}}^{(1)}}^{[r]}$ for the r th replicate, then the 95% confidence interval for $\hat{\lambda}_{g^{(1)}}^{(1)[r]}$ is con-

structured as $\text{CI}_{\lambda_{g^{(1)}}}^{[r]} = (\widehat{\lambda}_{g^{(1)}}^{(1)[r]} - 1.96 \times \widehat{\text{SE}}_{\lambda_{g^{(1)}}}^{[r]}, \widehat{\lambda}_{g^{(1)}}^{(1)[r]} + 1.96 \times \widehat{\text{SE}}_{\lambda_{g^{(1)}}}^{[r]})$. Here $\widehat{\text{SE}}_{\lambda_{g^{(1)}}}^{[r]}$ is obtained by Theorem 7, plugging in the variance estimator $\widehat{\sigma}^2$, which can be obtained by the procedure in Appendix B.2. Subsequently, the coverage probability (CP) is formed as $\text{CP}_{\lambda_{g^{(1)}}} = R^{-1} \sum_{r=1}^R I(\lambda_{g^{(1)}}^{(1)0} \in \text{CI}_{\lambda_{g^{(1)}}}^{[r]})$. We calculate the CPs for other parameters similarly. For comparison, we also calculate the RMSE and CP values for the oracle estimators under the true group memberships (denoted as $\widehat{\lambda}_{g^{(1)}}^{(1)\text{or}}, \widehat{\lambda}_{g^{(2)}}^{(2)\text{or}}, \widehat{\zeta}_{g^{(1)}}^{(1)\text{or}}, \widehat{\zeta}_{g^{(2)}}^{(2)\text{or}}, \widehat{\alpha}_{g^{(1)g^{(2)}}}^{\text{or}}$ accordingly). Lastly, to evaluate the group memberships estimation, we calculate the mis-clustering rates for two network groups as $\widehat{\eta}_l = (N_l R)^{-1} \sum_{r=1}^R \sum_{i_l} I(\widehat{g}_{i_l}^{(l)} \neq g_{i_l}^{(l)0})$ for $l = 1, 2$, where $\widehat{g}_{i_l}^{(l)}$ is the estimated group membership of the i_l th node. Here the mis-clustering rates are calculated after proper group permutations.

The simulation results of $G_{1,0} = G_{2,0} = 3$ are shown in Table 2. The first finding across all combinations is that once the group numbers are specified as the true values in advance, our iterative method can estimate the true group memberships with high accuracy, especially when the sample size is large. As N_1, N_2 or T increase, the mis-clustering rates for both network groups approach zero. Furthermore, we note that the RMSEs decrease either when the network sizes N_1 and N_2 increase or the time length T increases, and they approach the oracle RMSEs when the sample sizes are large. Next, we inspect the statistical inference results. We observe that in Table 2, the CPs are slightly smaller when the sample sizes are not very large, but they grow up to around 0.95 as N_1, N_2 and T increase. This guarantees that even under the scenario with a large number of parameters to be estimated, GTNAR can still perform well in terms of both estimation and inference for sufficiently large sample sizes. The patterns are similar in the other two scenarios.

6.2.2 ESTIMATION WHEN $G_l \geq G_{l,0}$ FOR $l = 1, 2$

We next consider the case of estimation without specifying the true group numbers in advance. Specifically, we estimate the group numbers by QIC in Section 3.3, where the tuning parameter is set to be $\kappa = 1/\{40\log(T)T^{1/8}\}$. Let the true group numbers be $G_{1,0} = G_{2,0} = 3$, and the corresponding true parameters are shown in Table 1. To evaluate the estimation accuracy, we calculate the RMSE for each parameter as explained below. Take $\boldsymbol{\lambda}^{(1)}$ for example, define $\text{RMSE}_{\boldsymbol{\lambda}^{(1)}, \text{all}} = \{(RN_1)^{-1} \sum_{r=1}^R \sum_{i_1=1}^{N_1} (\widehat{\lambda}_{g_{i_1}^{(1)}}^{(1)[r]} - \lambda_{g_{i_1}^{(1)0}}^{(1)0})^2\}^{1/2}$ as the RMSE for all nodes. RMSE for other parameters is calculated similarly. For the group memberships, the mis-clustering rates are calculated following the idea of Zhu et al. (2023). Recall that we have $\widehat{\mathcal{R}}_{g^{(l)}}^{(l)} = \{i_l : \widehat{g}_{i_l}^{(l)} = g^{(l)}\}$, where $\widehat{g}_{i_l}^{(l)}$ is denoted as the estimated group membership for the i_l th in the l th network. Note that G_l is not necessarily equal to $G_{l,0}$. In this case, we define the mappings from the estimated group memberships to the true group memberships $\chi_l : \{1, \dots, G_l\} \rightarrow \{1, \dots, G_{l,0}\}$ as $\chi_l(g^{(l)}) = \text{argmax}_{g^{(l)'} \in \{1, \dots, G_{l,0}\}} \sum_{i_l=1}^{N_l} I(i_l \in \widehat{\mathcal{R}}_{g^{(l)}}^{(l)}, g_{i_l}^{(l)0} = g^{(l)'})$ for $g^{(l)} \in \{1, \dots, G_l\}$. Thus, the mapping $\chi_l(g^{(l)})$ maps group $g^{(l)}$ to the true membership $g^{(l)'}$ where the majority of nodes in $\widehat{\mathcal{R}}_{g^{(l)}}^{(l)}$ belong to. Then, for the row group memberships, the mis-clustering rate in the r th replicate is defined as $\widehat{\xi}_l^{[r]} = N_l^{-1} \sum_{g^{(l)}=1}^{G_l} \sum_{i_l=1}^{N_l} I(i_l \in \widehat{\mathcal{R}}_{g^{(l)}}^{(l)[r]}, g_{i_l}^{(l)0} \neq \chi_l(g^{(l)}))$, where

Table 2: RMSEs ($\times 1000$) of estimated parameters under scenario 3 ($G_{1,0} = G_{2,0} = 3$) with 500 replications. The performances are evaluated for different sample sizes N_1, N_2 and the time length T . Results under two network structures are provided. The corresponding CPs are shown in parentheses.

Network	G_1	G_2	N_1	N_2	T	$\hat{\lambda}^{(1)}$	$\hat{\lambda}^{(2)}$	$\hat{\zeta}^{(1)}$	$\hat{\zeta}^{(2)}$	$\hat{\alpha}$	$\hat{\lambda}^{(1)\text{or}}$	$\hat{\lambda}^{(2)\text{or}}$	$\hat{\zeta}^{(1)\text{or}}$	$\hat{\zeta}^{(2)\text{or}}$	$\hat{\alpha}^{\text{or}}$	$\hat{\eta}_1$	$\hat{\eta}_2$
SBM			100	80	20	25.6 (0.705)	12.5 (0.907)	59.1 (0.595)	15.9 (0.920)	58.6 (0.731)	11.7 (0.947)	10.8 (0.938)	13.6 (0.945)	12.9 (0.949)	19.1 (0.949)	0.1634	0.0113
					40	9.1 (0.911)	7.5 (0.945)	14.4 (0.871)	9.1 (0.948)	15.2 (0.913)	7.8 (0.951)	7.4 (0.945)	9.4 (0.938)	9.1 (0.947)	0.0285	0.0001	
		3	200	150	20	9.1 (0.825)	5.6 (0.925)	17.4 (0.805)	7.0 (0.940)	14.5 (0.857)	6.1 (0.932)	5.3 (0.939)	6.8 (0.940)	6.9 (0.942)	9.6 (0.945)	0.0524	0.0006
					40	4.2 (0.942)	3.8 (0.939)	5.4 (0.940)	4.7 (0.949)	7.0 (0.941)	4.0 (0.948)	3.8 (0.941)	4.6 (0.950)	6.6 (0.948)	0.0033	0.0001	
			300	250	20	4.4 (0.916)	3.5 (0.935)	6.1 (0.906)	4.2 (0.941)	6.9 (0.920)	3.7 (0.948)	3.5 (0.936)	4.2 (0.946)	4.2 (0.941)	5.9 (0.944)	0.0092	0.0001
					40	3.0 (0.930)	2.3 (0.948)	4.2 (0.927)	2.9 (0.952)	4.8 (0.936)	2.5 (0.947)	2.3 (0.950)	2.9 (0.948)	4.1 (0.948)	0.0006	0	
Power-Law			100	80	20	17.5 (0.769)	9.2 (0.931)	40.2 (0.687)	17.6 (0.923)	34.7 (0.790)	11.0 (0.938)	8.6 (0.941)	12.9 (0.948)	13.2 (0.944)	17.8 (0.947)	0.1361	0.0156
					40	11.4 (0.863)	5.7 (0.947)	20.0 (0.831)	9.3 (0.937)	19.5 (0.886)	7.2 (0.945)	5.6 (0.943)	9.3 (0.941)	13.2 (0.940)	0.0411	0	
		3	200	150	20	7.7 (0.843)	4.6 (0.939)	14.1 (0.840)	6.8 (0.945)	12.1 (0.883)	5.4 (0.937)	4.5 (0.945)	6.7 (0.947)	6.7 (0.946)	8.7 (0.950)	0.0408	0.0001
					40	4.3 (0.933)	2.1 (0.948)	5.7 (0.927)	4.7 (0.948)	6.6 (0.940)	4.0 (0.946)	2.1 (0.949)	4.7 (0.940)	6.2 (0.948)	0.0039	0	
			300	250	20	4.7 (0.890)	2.4 (0.937)	7.8 (0.892)	5.4 (0.938)	10.7 (0.907)	3.3 (0.939)	2.1 (0.943)	4.2 (0.943)	4.3 (0.944)	5.5 (0.945)	0.0170	0.0001
					40	2.7 (0.942)	1.9 (0.936)	3.2 (0.950)	2.9 (0.948)	4.1 (0.947)	2.6 (0.944)	1.9 (0.935)	2.9 (0.948)	4.0 (0.949)	0.0010	0	

$\widehat{\mathcal{R}}_{g^{(l)}}^{(l)[r]}$ is the estimated node set belonging to the group $g^{(l)}$ in the r th replicate. Then, the overall group memberships error rate is calculated as $\widehat{\xi}_l = R^{-1} \sum_r \widehat{\xi}_l^{[r]}$. The results are shown in Table 3.

We discuss the results shown in Table 3 from two aspects. On one hand, when the group number is under-specified ($G_1 = 2, G_2 = 2$), the node-wise RMSEs are large and usually do not decrease when the N_1, N_2 and T grow. Besides, the error rates $\widehat{\xi}_1$ and $\widehat{\xi}_2$ are around 0.3, indicating a low accuracy in estimating node memberships. These results are expected since a non-ignorable estimation bias exists in an under-fitted model. On the other hand, when the group numbers G_l s are correctly ($G_1 = 3, G_2 = 3$) or over-specified ($G_1 = 4, G_2 = 4$), the RMSE values are generally much lower. This is consistent with our theoretical analysis in Theorem 2.

7. Real Data Applications

7.1 Data Description

The Yelp dataset spans from 2010 to 2018 and covers five North American cities: Charlotte, Las Vegas, Phoenix, Scottsdale, and Toronto. The observation period is divided into $T = 36$ quarters. To ensure data quality, we filter the dataset to retain active users who have provided more than 5 reviews over this time span. We further divide each city into districts, as illustrated in Figure 3(a). Our response variable, denoted as $Y_{ij,t}$, represents the $\log(1+x)$ transformed number of reviews by user i on district j during the t th quarter. Here $Y_{ij,t}$ is treated as a continuous variable in our real data analysis. To visualize temporal trends in review activity, we calculate the quarterly average responses for each city and depict them in Figure 3(b). Different patterns emerge from this analysis. For instance, Las Vegas stands out as the city with the most reviews, reflecting its bustling business environment. Charlotte, Phoenix, and Scottsdale exhibit relatively similar and stable review trends. In contrast, Toronto shows a noticeable increase in review volume after 2015, likely due to Yelp’s expansion in the Toronto area during that period.

Next, we construct the adjacency matrices for users ($\mathbf{A}^{(1)}$) and districts ($\mathbf{A}^{(2)}$) as follows. The user network is built based on the friend list information. Specifically, if user j is on the friend list of user i on Yelp, then we set $a_{ij}^{(1)} = 1$. Otherwise we set $a_{ij}^{(1)} = 0$. The spatial network is built based on the geographical adjacent relationship. Specifically, we set $a_{ij}^{(2)} = 1$ if the district j is adjacent to district i .

Lastly, to characterize the dynamic patterns of the responses, we collect a number of covariates for users and districts, respectively. For user i in quarter t , we consider the following five covariates: (1) the number of months after joining Yelp by the start of the quarter t ($x_{it,\text{dur}}^{(1)}$), (2) whether the user is VIP by the start of the quarter t ($x_{it,\text{vip}}^{(1)}$), (3) average tags (i.e., “useful”, “funny” and “cool”) the user i obtains for his/her reviews during the last quarter ($x_{it,\text{use}}^{(1)}, x_{it,\text{fun}}^{(1)}, x_{it,\text{cool}}^{(1)}$). Next, for the j th district in quarter t , we consider two covariates: (1) the average “stars” ($x_{jt,\text{star}}^{(2)}$), and (2) the average review number ($x_{jt,\text{num}}^{(2)}$) obtained by the j th district during the $(t-1)$ th quarter. These two covariates are indicative of the average popularity levels in the preceding time period. We standardize all continuous covariates to the range $[0, 1]$ for subsequent analysis. In Figure A.11 in Appendix M.1, we

Table 3: Simulation results for the two network examples with pre-specified group numbers as well as the QIC selection group numbers \widehat{G}_l . The true group numbers are set as $G_{1,0} = G_{2,0} = 3$. The node-wise RMSEs of different estimators are denoted as $\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(1)}$, $\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(2)}$, $\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(1)}$, $\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(2)}$, $\widehat{\boldsymbol{\alpha}}_{\text{all}}$.

N_1	N_2	T	G_1	G_2	Scenario 1 (SBM)					Scenario 2 (Power-Law)									
					$\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(1)}$	$\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(2)}$	$\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(1)}$	$\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(2)}$	$\widehat{\boldsymbol{\alpha}}_{\text{all}}$	$\widehat{\xi}_1$	$\widehat{\xi}_2$	$\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(1)}$	$\widehat{\boldsymbol{\lambda}}_{\text{all}}^{(2)}$	$\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(1)}$	$\widehat{\boldsymbol{\zeta}}_{\text{all}}^{(2)}$	$\widehat{\boldsymbol{\alpha}}_{\text{all}}$	$\widehat{\xi}_1$	$\widehat{\xi}_2$	
100	80	20	Oracle		0.0065	0.0059	0.0071	0.0070	0.0061	-	-	0.0072	0.0039	0.0071	0.0072	0.0060	-	-	
			2	2	0.0285	0.0419	0.0427	0.0583	0.0368	0.3219	0.3125	0.0237	0.0438	0.0382	0.0647	0.0367	0.2705	0.3705	
			3	3	0.0207	0.0096	0.0251	0.0099	0.0174	0.1403	0.0142	0.0196	0.0129	0.0263	0.0114	0.0169	0.1487	0.0208	
			4	4	0.0192	0.0099	0.0181	0.0097	0.0172	0.0720	0.0042	0.0201	0.0134	0.0223	0.0101	0.0176	0.1046	0.0077	
	\widehat{G}_1	\widehat{G}_2	0.0278	0.0407	0.0416	0.0565	0.0358	-	-	0.0218	0.0374	0.0339	0.0554	0.0323	-	-			
	40	80	40	Oracle		0.0044	0.0043	0.0050	0.0050	0.0043	-	-	0.0046	0.0040	0.0050	0.0049	0.0040	-	-
				2	2	0.0216	0.0436	0.0430	0.0503	0.0327	0.3100	0.2625	0.0210	0.0428	0.0403	0.0607	0.0361	0.2800	0.3375
				3	3	0.0082	0.0070	0.0104	0.0063	0.0078	0.0366	0.0061	0.0069	0.0040	0.0085	0.0049	0.0053	0.0233	0
				4	4	0.0081	0.0061	0.0084	0.0065	0.0084	0.0136	0.0006	0.0076	0.0066	0.0078	0.0067	0.0082	0.0101	0.0015
				\widehat{G}_1	\widehat{G}_2	0.0086	0.0112	0.0121	0.0121	0.0100	-	-	0.0069	0.0040	0.0085	0.0049	0.0053	-	-
Oracle					0.0034	0.0031	0.0036	0.0037	0.0031	-	-	0.0036	0.0021	0.0036	0.0037	0.0030	-	-	
200	150	20	2	2	0.0221	0.0431	0.0455	0.0592	0.0351	0.3505	0.3133	0.0217	0.0420	0.0416	0.0559	0.0351	0.3002	0.3200	
			3	3	0.0091	0.0037	0.0112	0.0040	0.0068	0.0504	0.0015	0.0076	0.0022	0.0089	0.0037	0.0054	0.0342	0.0001	
			4	4	0.0106	0.0052	0.0089	0.0054	0.0096	0.0296	0.0003	0.0086	0.0047	0.0082	0.0050	0.0077	0.0240	0.0004	
			\widehat{G}_1	\widehat{G}_2	0.0178	0.0334	0.0354	0.0457	0.0275	-	-	0.0076	0.0026	0.0089	0.0042	0.0056	-	-	
	40	150	40	Oracle		0.0024	0.0021	0.0026	0.0026	0.0022	-	-	0.0023	0.0011	0.0025	0.0026	0.0020	-	-
				2	2	0.0199	0.0405	0.0383	0.0539	0.0343	0.2600	0.3133	0.0218	0.0409	0.0437	0.0554	0.0350	0.3250	0.3133
				3	3	0.0029	0.0021	0.0034	0.0026	0.0026	0.0028	0.0001	0.0026	0.0011	0.0030	0.0026	0.0022	0.0027	0
				4	4	0.0044	0.0038	0.0043	0.0040	0.0055	0.0047	0.0017	0.0049	0.0042	0.0045	0.0042	0.0057	0.0063	0.0032
				\widehat{G}_1	\widehat{G}_2	0.0030	0.0027	0.0035	0.0034	0.0029	-	-	0.0026	0.0011	0.0030	0.0026	0.0022	-	-

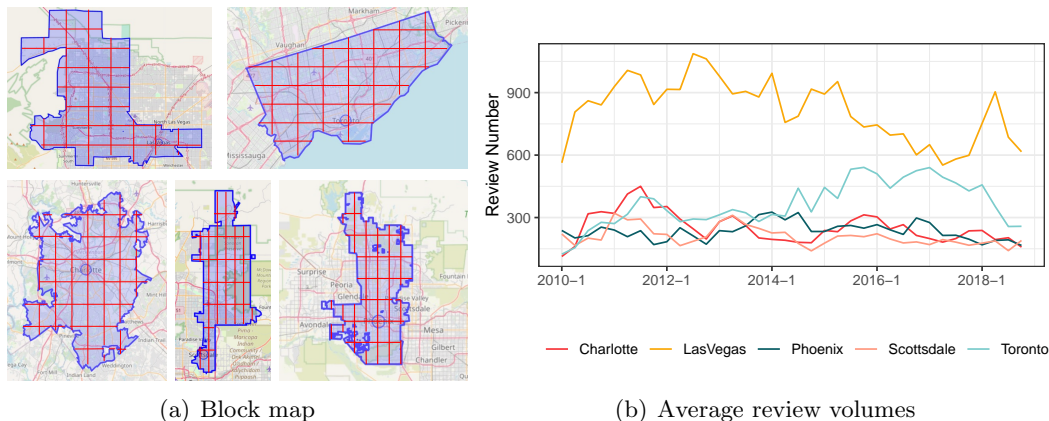


Figure 3: (a) Geographical maps with split districts in each city. The city from the top left panel to the bottom right panel shows the map of Las Vegas, Toronto, Charlotte, Scottsdale, and Phoenix, respectively. (b) Average number of reviews from 2010-Q1 to 2018-Q4 in five cities.

visualize the relationship between these user-related covariates and the response variable. The plot reveals that users who receive more tags for their reviews tend to be motivated to contribute more reviews. Notably, VIP users in Scottsdale and Toronto tend to write more reviews, whereas VIP users in Charlotte exhibit comparatively less activity. Subsequently, we apply the GTNAR model to each of the five cities, enabling us to analyze and understand the distinctive group patterns within each urban area.

7.2 Estimation Results

There has been a lot of research finding that the activeness of users on the platform can lead to higher profit (Forman et al., 2008; Pansari and Kumar, 2017), hence finding out the key factors that positively affect the review volumes can help business owners make personalized recommendations to users and develop a commercial strategy. We employ QIC for the selection of group numbers, and show the results of Charlotte, Las Vegas, and Phoenix in Table 4. The results of Scottsdale and Toronto are provided in Table A.16 in Appendix M.1. The numbers of user groups and district groups vary across the five cities, indicating different levels of heterogeneity among them. For instance, consider the results for Phoenix, where there are estimated 3 user groups and 2 district groups.

Notably, the spatial (column) network effects are consistently positive, suggesting a favorable effect from neighboring districts. Such an observation is consistent with the findings in the literature (e.g. Sun and Paule, 2017). Furthermore, we can also observe that within the two group-wise spatial effects, $\hat{\lambda}_2^{(2)} = 0.270$ is larger than that of Group 1 ($\hat{\lambda}_1^{(2)} = 0.05$), signifying a stronger neighbor effect. In other words, if the second group of districts' neighbors obtain more reviews in the last period, then it is likely that these districts would receive more reviews in this period. We further visualize the districts by estimated groups on the left panel of Figure 4, where the gray districts are from Group 2, and the white districts are from Group 1. Subsequently, we mark the shops on the map, which are shown by red

points in the right panel of Figure 4. By scrutinizing the locations, we find that the shops in Group 2 are mainly located in the central business districts. The average number of shops per district in Group 2 is 25.16, whereas that in Group 1 is 6.57, showing a higher density of shops in the gray areas (Group 2). Our estimation results reveal that shops in Group 2 exhibit larger spillover effects on their neighbors. This finding is consistent with the existing research, which also reveals a stronger spillover effect for geometrically clustered units (Arzaghi and Henderson, 2008; Rossi-Hansberg et al., 2010; Vitorino, 2012). Performing a

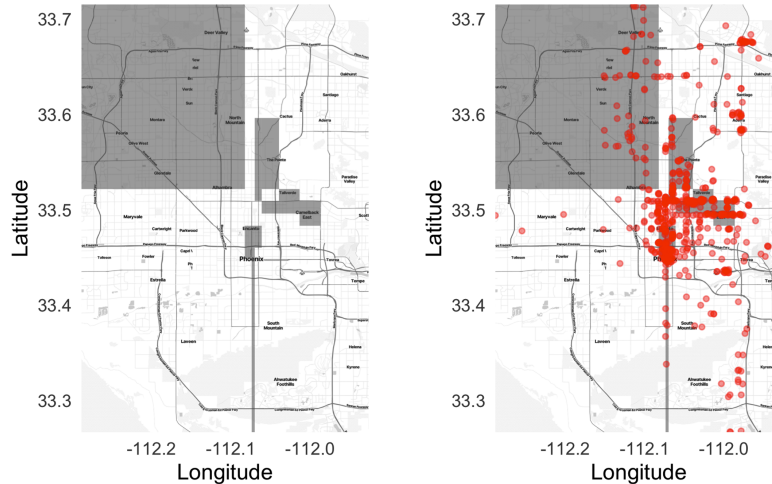


Figure 4: The left panel shows the two groups of districts in Phoenix (Group 2 is marked as gray). The right panel marks all shops as red points.

similar investigation on the user dimension, one can find that the social network effect in Group 1 appears to be significantly negative ($\hat{\lambda}_1^{(1)} = -0.02$), whereas corresponding coefficients in the other two groups are positive. This implies that user activities in Group 1 are influenced oppositely by their friends' behaviors, whereas in the other groups, users are still positively influenced by their friends. By scrutinizing deeper into users in Group 1, we found that their average duration after joining Yelp is longer than in Groups 2 and 3. This finding aligns with existing works showing that the social effects would decrease as the users' membership duration prolongs (Nitzan and Libai, 2011; Aral and Walker, 2011). Further, the social network effect is the largest in Group 3, i.e., $\hat{\lambda}_3^{(1)} = 0.024$. The network in-degree and out-degree of users in this group are higher than in Groups 1 and 2, which is visualized in the left panel of Figure 5. The phenomenon of larger network degrees associated with larger social network effects has been extensively validated in literature (Hill et al., 2006; Hinz et al., 2011; Susarla et al., 2012). Therefore, the business owners can implement a precision marketing strategy targeting users in Group 3, who have stronger network effects.

Besides, the estimated $\hat{\alpha}$ values are all positive, indicating a positive self-motivated effect overall, whose values are visualized in the right panel of Figure 5. Specifically, the momentum effect on the review number is the largest for the Group 1 users' evaluation on

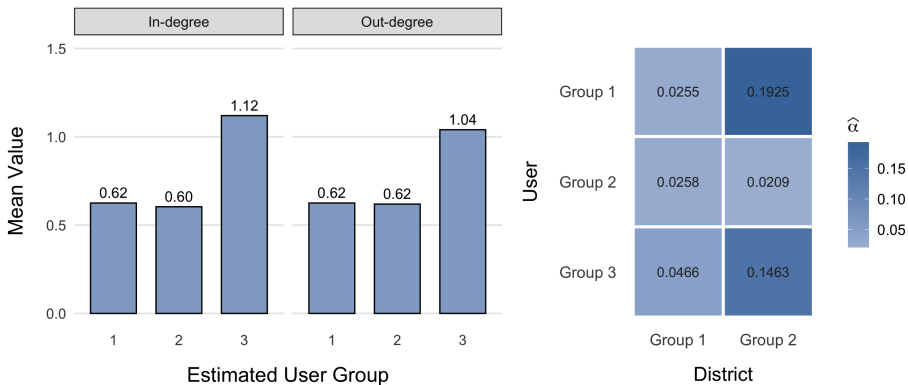


Figure 5: Left panel: average in-degree and out-degree of users in the three estimated groups in Phoenix. Right panel: heatmap of estimated momentum effect $\hat{\alpha}$ in Phoenix.

the Group 2 districts, i.e., $\hat{\alpha}_{12} = 0.193$. This implies that users in Group 1 have higher self-driven persistency in the historical reviews than the other groups, while they exhibit smaller social network effects.

8. Concluding Remarks

In this work, we introduce a novel Group Tensor Network Autoregression (GTNAR) model designed for time series data indexed by multi-relational networks. By leveraging network structures on each dimension, GTNAR establishes a unique framework for analyzing tensor-valued time series data. This model presents a valuable tool for analyzing data collected in complex network environments, shedding light on various network effects. From a modeling perspective, several intriguing future topics emerge. It would be interesting to consider the high-dimensional covariates in this framework. Furthermore, it is worth exploring models that consider multiple network effects in a multiplicative form, as investigated by Chen et al. (2021). Besides, introducing a hidden factor structure into the GTNAR model can potentially offer more insights into high-dimensional data, capturing more underlying information within the tensor-valued time series, and thus represents an interesting future research topic. Lastly, it is also important to further investigate how to model categorical responses in our modeling framework. From the theoretical perspective, although the convergence rate established in Theorem 2 is typical in existing research, how to obtain a minimax rate needs further investigation.

Acknowledgments

Yimeng Ren and Xuening Zhu’s research are supported by the National Natural Science Foundation of China (nos. 72573038, 12331009), MOE Laboratory for National Devel-

opment and Intelligent Governance, Fudan University. Yanyuan Ma's work is partially supported by grants from NIH.

References

- Radoslaw Adamczak. A note on the hanson-wright inequality for random vectors with dependencies. *Electronic Communications in Probability*, 20:1–13, 2015.
- Gediminas Adomavicius, Ramesh Sankaranarayanan, Shahana Sen, and Alexander Tuzhilin. Incorporating contextual information in recommender systems using a multidimensional approach. *ACM Transactions on Information systems (TOIS)*, 23(1):103–145, 2005.
- Luiz GA Alves, Giuseppe Mangioni, Isabella Cingolani, Francisco Aparecido Rodrigues, Pietro Panzarasa, and Yamir Moreno. The nested structural organization of the worldwide trade multi-layer network. *Scientific reports*, 9(1):2866, 2019.
- Tomohiro Ando and Jushan Bai. Panel data models with grouped factor structure under unknown group membership. *Journal of Applied Econometrics*, 31(1):163–191, 2016.
- Tomohiro Ando and Jushan Bai. Quantile co-movement in financial markets: A panel quantile model with unobserved heterogeneity. *Journal of the American Statistical Association*, 115(529):266–279, 2020.
- Sinan Aral and Dylan Walker. Creating social contagion through viral product design: A randomized trial of peer influence in networks. *Management science*, 57(9):1623–1639, 2011.
- Mirko Armillotta and Konstantinos Fokianos. Nonlinear network autoregression. *The Annals of Statistics*, 51(6):2526–2552, 2023.
- Mohammad Arzaghi and J Vernon Henderson. Networking off madison avenue. *The Review of Economic Studies*, 75(4):1011–1038, 2008.
- Marc Auboin and Michele Ruta. The relationship between exchange rates and international trade: a literature review. *World Trade Review*, 12(3):577–605, 2013.
- Sumanta Basu, George Michailidis, et al. Regularized estimation in sparse high-dimensional time series models. *The Annals of Statistics*, 43(4):1535–1567, 2015.
- Nicolas Berman and Jérôme Héricourt. Financial factors and the margins of trade: Evidence from cross-country firm-level data. *Journal of Development Economics*, 93(2):206–217, 2010.
- Andrew B Bernard and J Bradford Jensen. Exceptional exporter performance: cause, effect, or both? *Journal of international economics*, 47(1):1–25, 1999.
- Andrew B Bernard and J Bradford Jensen. Why some firms export. *Review of economics and Statistics*, 86(2):561–569, 2004.
- C Alan Bester and Christian B Hansen. Grouped effects estimators in fixed effects models. *Journal of Econometrics*, 190(1):197–208, 2016.
- Stéphane Bonhomme and Elena Manresa. Grouped patterns of heterogeneity in panel data. *Econometrica*, 83(3):1147–1184, 2015.

- Raymond J Carroll and Daren BH Cline. An asymptotic theory for weighted least-squares with weights estimated by replication. *Biometrika*, 75(1):35–43, 1988.
- Elias Carroni, Paolo Pin, and Simone Righi. Bring a friend! privately or publicly? *Management Science*, 66(5):2269–2290, 2020.
- Jinyuan Chang, Jing He, Lin Yang, and Qiwei Yao. Modelling matrix time series via a tensor cp-decomposition. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(1):127–148, 2023.
- Elynn Y Chen and Jianqing Fan. Statistical inference for high-dimensional matrix-variate factor models. *Journal of the American Statistical Association*, 118(542):1038–1055, 2023.
- Elynn Y Chen, Ruey S Tsay, and Rong Chen. Constrained factor models for high-dimensional matrix-variate time series. *Journal of the American Statistical Association*, 115(530):775–793, 2020a.
- Elynn Y Chen, Ruey S Tsay, and Rong Chen. Constrained factor models for high-dimensional matrix-variate time series. *Journal of the American Statistical Association*, 2020b.
- Elynn Y Chen, Jianqing Fan, and Xuening Zhu. Community network auto-regression for high-dimensional time series. *Journal of Econometrics*, 235(2):1239–1256, 2023.
- Jian Chen, Weidong Ma, and Ganggang Xu. Group network multivariate garch. *Available at SSRN 6760338*, 2026.
- Rong Chen, Han Xiao, and Dan Yang. Autoregressive models for matrix-valued time series. *Journal of Econometrics*, 222(1):539–560, 2021.
- Rong Chen, Dan Yang, and Cun-Hui Zhang. Factor models for high-dimensional tensor time series. *Journal of the American Statistical Association*, 117(537):94–116, 2022.
- Weilin Chen and Clifford Lam. Rank and factor loadings estimation in time series tensor factor model by pre-averaging. *The Annals of Statistics*, 52(1):364–391, 2024.
- Xi Chen, Ralf Van Der Lans, and Tuan Q Phan. Uncovering the importance of relationship characteristics in social networks: Implications for seeding strategies. *Journal of Marketing Research*, 54(2):187–201, 2017.
- Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. Power-law distributions in empirical data. *SIAM Review*, 51:661–703, 2009.
- Enrico Corradini, Antonino Nocera, Domenico Ursino, and Luca Virgili. Investigating negative reviews and detecting negative influencers in yelp through a multi-dimensional social network based model. *International Journal of Information Management*, 60:102377, 2021.
- Richard A Davis, Pengfei Zang, and Tian Zheng. Sparse vector autoregressive modeling. *Journal of Computational and Graphical Statistics*, 25(4):1077–1096, 2016.

- Shanshan Ding and R Dennis Cook. Matrix variate regressions and envelope models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 80(2):387–408, 2018.
- Michael Emch, Elisabeth D Root, Sophia Giebultowicz, Mohammad Ali, Carolina Perez-Heydrich, and Mohammad Yunus. Integration of spatial and social network analysis in disease transmission studies. In *Geographies of Health, Disease and Well-being*, pages 130–141. Routledge, 2016.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- Jianqing Fan, Yuan Ke, and Yuan Liao. Augmented factor models with applications to validating market risk factors and forecasting bond risk premia. *Journal of Econometrics*, 222(1):269–294, 2021.
- Guanhua Fang, Ganggang Xu, Haochen Xu, Xuening Zhu, and Yongtao Guan. Group network hawkes process. *Journal of the American Statistical Association*, pages 1–17, 2023.
- Hao Fe. Social networks and consumer behavior: evidence from yelp. *Journal of Economic Behavior & Organization*, 209:1–14, 2023.
- Yusen Feng, Gang-Jin Wang, You Zhu, and Chi Xie. Systemic risk spillovers and the determinants in the stock markets of the belt and road countries. *Emerging Markets Review*, 55:101020, 2023.
- Denzil G Fiebig, Michael P Keane, Jordan Louviere, and Nada Wasi. The generalized multinomial logit model: accounting for scale and coefficient heterogeneity. *Marketing science*, 29(3):393–421, 2010.
- Chris Forman, Anindya Ghose, and Batia Wiesenfeld. Examining the relationship between reviews and sales: The role of reviewer identity disclosure in electronic markets. *Information Systems Research*, 19(3):291–313, 2008.
- Caroline L Freund and Diana Weinhold. The effect of the internet on international trade. *Journal of international economics*, 62(1):171–189, 2004.
- Chang Gan, Mihai Voda, Kai Wang, Lijun Chen, and Jun Ye. Spatial network structure of the tourism economy in urban agglomeration: A social network analysis. *Journal of Hospitality and Tourism Management*, 47:124–133, 2021.
- Peter Hall and Christopher C Heyde. *Martingale limit theory and its application*. Academic press, 2014.
- Yuefeng Han, Rong Chen, Cun-Hui Zhang, and Qiwei Yao. Simultaneous decorrelation of matrix time series. *Journal of the American Statistical Association*, pages 1–13, 2023.
- Yuefeng Han, Dan Yang, Cun-Hui Zhang, and Rong Chen. Cp factor model for dynamic tensors. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkae036, 2024.

- David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3): 1079–1083, 1971.
- Shawndra Hill, Foster Provost, and Chris Volinsky. Network-Based Marketing: Identifying Likely Adopters via Consumer Networks. *Statistical Science*, 21(2):256–276, 2006.
- Oliver Hinz, Bernd Skiera, Christian Barrot, and Jan U Becker. Seeding strategies for viral marketing: an empirical comparison. *Journal of Marketing*, 75(6):55–71, 2011.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Peter D Hoff. Multilinear tensor regression for longitudinal relational data. *The annals of applied statistics*, 9(3):1169, 2015.
- Peter Hooper and Steven W Kohlhagen. The effect of exchange rate uncertainty on the prices and volume of international trade. *Journal of international Economics*, 8(4):483–511, 1978.
- IMF. *Direction of Trade Statistics, International Monetary Fund*, 2017.
- Srikanth Jagabathula, Lakshminarayanan Subramanian, and Ashwin Venkataraman. A model-based embedding technique for segmenting customers. *Operations Research*, 66(5):1247–1267, 2018.
- Bing-Yi Jing, Ting Li, Zhongyuan Lyu, and Dong Xia. Community detection on mixture multilayer networks via regularized tensor decomposition. *The Annals of Statistics*, 49(6):3181–3205, 2021.
- Zheng Tracy Ke, Jianqing Fan, and Yichao Wu. Homogeneity pursuit. *Journal of the American Statistical Association*, 110(509):175–194, 2015.
- Julia Koschinsky. Spatial heterogeneity in spillover effects of assisted and unassisted rental housing. *Journal of Urban Affairs*, 31(3):319–347, 2009.
- Jing Lei, Kehui Chen, and Brian Lynch. Consistent community detection in multi-layer network data. *Biometrika*, 107(1):61–73, 2020.
- Chenlei Leng and Cheng Yong Tang. Sparse matrix graphical models. *Journal of the American Statistical Association*, 107(499):1187–1200, 2012.
- Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on learning theory*, pages 2738–2779. PMLR, 2020.
- Ruiqi Liu, Zuofeng Shang, Yonghui Zhang, and Qiankun Zhou. Identification and estimation in panel models with overspecified number of groups. *Journal of Econometrics*, 215(2): 574–590, 2020.
- Wenyang Liu, Ganggang Xu, Jianqing Fan, and Xuening Zhu. Two-way homogeneity pursuit for quantile network vector autoregression. *arXiv preprint arXiv:2404.18732*, 2024.

- Gábor Lugosi and Shahar Mendelson. Sub-gaussian estimators of the mean of a random vector. *The Annals of Statistics*, 47(2):783–794, 2019.
- Helmut Lütkepohl. *New introduction to multiple time series analysis*. Springer Science & Business Media, 2005.
- Zongming Ma and Sagnik Nandy. Community detection with contextual multilayer networks. *IEEE Transactions on Information Theory*, 69(5):3203–3239, 2023.
- Peter W MacDonald, Elizaveta Levina, and Ji Zhu. Latent space models for multiplex networks with shared structure. *Biometrika*, 109(3):683–706, 2022.
- Ke Miao, Peter CB Phillips, and Liangjun Su. High-dimensional vars with common factors. *Journal of Econometrics*, 233(1):155–183, 2023.
- Katta G. Murty and Santosh N. Kabadi. Some np-complete problems in quadratic and nonlinear programming. *Math. Program.*, 39(2):117–129, June 1987. ISSN 0025-5610.
- Mark Newman, Albert-Laszlo Barabasi, and Duncan J Watts. *The structure and dynamics of networks*. Princeton University Press, 2006.
- William B Nicholson, Ines Wilms, Jacob Bien, and David S Matteson. High dimensional forecasting via interpretable vector autoregression. *Journal of Machine Learning Research*, 21(166):1–52, 2020.
- Irit Nitzan and Barak Libai. Social effects on customer retention. *Journal of Marketing*, 75(6):24–38, 2011.
- Krzysztof Nowicki and Tom A B Snijders. Estimation and prediction for stochastic block-structures. *Journal of the American Statistical Association*, 96(455):1077–1087, 2001.
- Anita Pansari and Vera Kumar. Customer engagement: the construct, antecedents, and consequences. *Journal of the academy of marketing science*, 45(3):294–311, 2017.
- Stephen Redding and Anthony J Venables. Economic geography and international inequality. *Journal of international Economics*, 62(1):53–82, 2004.
- Esteban Rossi-Hansberg, Pierre-Daniel Sarte, and Raymond Owens III. Housing externalities. *Journal of political Economy*, 118(3):485–535, 2010.
- Paul-Marie Samson. Concentration of measure inequalities for markov chains and ϕ -mixing processes. *The Annals of Probability*, 28(1):416–461, 2000.
- Steven G Self and Kung-Yee Liang. Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association*, 82(398):605–610, 1987.
- Juan Camilo Serpa and Harish Krishnan. The impact of supply chains on firm-level productivity. *Management Science*, 64(2):511–532, 2018.

- Jun Shao. Asymptotic distribution of the weighted least squares estimator. *Annals of the Institute of Statistical Mathematics*, 41:365–382, 1989.
- Liangjun Su, Zhentao Shi, and Peter CB Phillips. Identifying latent structures in panel data. *Econometrica*, 84(6):2215–2264, 2016.
- Liangjun Su, Wuyi Wang, and Xingbai Xu. Identifying latent group structures in spatial dynamic panels. *Journal of Econometrics*, 235(2):1955–1980, 2023.
- Yeran Sun and Jorge David Gonzalez Paule. Spatial analysis of users-generated ratings of yelp venues. *Open Geospatial Data, Software and Standards*, 2:1–9, 2017.
- Anjana Susarla, Jeong-Ha Oh, and Yong Tan. Social networks and the diffusion of user-generated content: Evidence from youtube. *Information Systems Research*, 23(1):23–41, 2012.
- Michel Talagrand. An isoperimetric theorem on the cube and the kintchine-kahane inequalities. *Proceedings of the American Mathematical Society*, 104(3):905–909, 1988.
- Nuray Terzi. The impact of e-commerce on international trade and employment. *Procedia-social and behavioral sciences*, 24:745–753, 2011.
- Ashutosh Tiwari and Timothy J Richards. Social networks and restaurant ratings. *Agribusiness*, 32(2):153–174, 2016.
- Maria Ana Vitorino. Empirical entry games with complementarities: an application to the shopping center industry. *Journal of Marketing Research*, 49(2):175–191, 2012.
- Michael Vogt and Oliver Linton. Classification of non-parametric regression functions in longitudinal data models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 79(1):5–27, 2017.
- AT Walden and A Serroukh. Wavelet analysis of matrix-valued time-series. *Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 458(2017):157–179, 2002.
- Di Wang and Ruey S Tsay. Rate-optimal robust estimation of high-dimensional vector autoregressive models. *The Annals of Statistics*, 51(2):846–877, 2023.
- Di Wang, Yao Zheng, Heng Lian, and Guodong Li. High-dimensional vector autoregressive time series modeling via tensor decomposition. *Journal of the American Statistical Association*, 117(539):1338–1356, 2022.
- Di Wang, Yao Zheng, and Guodong Li. High-dimensional low-rank tensor autoregressive time series modeling. *Journal of Econometrics*, 238(1):105544, 2024.
- Dong Wang, Xialu Liu, and Rong Chen. Factor models for matrix-valued high-dimensional time series. *Journal of Econometrics*, 208(1):231–248, 2019.
- Lan Wang, Yongdai Kim, and Runze Li. Calibrating non-convex penalized regression in ultra-high dimension. *The Annals of Statistics*, 41(5):2505–2536, 2013.

- Michel Wedel and Wagner A Kamakura. *Market segmentation: Conceptual and methodological foundations*. Springer Science & Business Media, 2000.
- Shu Wu, Qiang Liu, Liang Wang, and Tieniu Tan. Contextual operation for recommender systems. *IEEE Transactions on Knowledge and Data Engineering*, 28(8):2000–2012, 2016.
- Zhonghua Yin, Fang Wang, and Jianbang Gan. Spatial spillover effects of global forest product trade. *Forest Policy and Economics*, 113:102112, 2020. ISSN 1389-9341.
- Xuefei Zhang, Songkai Xue, and Ji Zhu. A flexible latent space model for multilayer networks. In *International Conference on Machine Learning*, pages 11288–11297. PMLR, 2020.
- Yingying Zhang, Huixia Judy Wang, and Zhongyi Zhu. Quantile-regression-based clustering for panel data. *Journal of Econometrics*, 213(1):54–67, 2019.
- Hua Zhou and Lexin Li. Regularized matrix regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(2):463–483, 2014.
- Jie Zhou, Will Wei Sun, Jingfei Zhang, and Lexin Li. Partially observed dynamic tensor response regression. *Journal of the American Statistical Association*, 118(541):424–439, 2023.
- Shuheng Zhou. Gemini: Graph estimation with matrix variate normal instances. *The Annals of Statistics*, 42(2):532–562, 2014.
- Xuening Zhu and Rui Pan. Grouped network vector autoregression. *Statistica Sinica*, 30(3):1437–1462, 2020.
- Xuening Zhu, Rui Pan, Guodong Li, Yuewen Liu, and Hansheng Wang. Network vector autoregression. *The Annals of Statistics*, 45(3):1096–1123, 2017.
- Xuening Zhu, Ganggang Xu, and Jianqing Fan. Simultaneous estimation and group identification for network vector autoregressive model with heterogeneous nodes. *Journal of Econometrics*, page 105564, 2023.
- Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429, 2006.

In this part, the required notations are provided in Appendix A. More estimation details for $q = 2$ and general q are provided in Appendix B–D. An numerical convergence analysis for the algorithm is shown in Appendix E. The initialization procedure and of our main algorithm is provided in Appendix F. We provide additional discussion about the interactive model in Appendix G, followed by discussion on future extensions in Appendix H. We provide the technical details and proofs for the main text in Section I, and several useful lemmas are given in Section J–K. Besides, we provide a number of additional simulation studies in Appendix L. Lastly, some additional real data analysis are given in Appendix M and Appendix N.

Appendix A. Notations

Define $Q^*(\Theta) = E\{Q(\Theta)\}$ and $Q_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q}) = E\{Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q})\}$, where $Q(\Theta)$ and $Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q})$ are as defined in (16) in the main text. Denote $\mathcal{G}_{-l} = \{\mathcal{G}_m : m \neq l\}$ and $\mathbf{i}_{-l} = (i_m : m \neq l)^\top \in \mathbb{R}^{q-1}$. Further, denote $\boldsymbol{\theta}^{-(l)} = \{(\boldsymbol{\theta}_{g^{(m)}}^{(m)} : g^{(m)} \in [G_m]) \in \mathbb{R}^{G_m(p_m+1)} : m \neq l, 1 \leq m \leq q\}$, $\boldsymbol{\xi}_{g_{i_l}}^{(l)} = (\boldsymbol{\theta}_{g_{i_l}}^{(l)\top}, \text{vec}(\boldsymbol{\alpha}_{g_{i_l}}))^\top$, $\boldsymbol{\xi}_{g^{-(l)}} = (\boldsymbol{\theta}_{g^{(m)}}^{(m)\top}, \text{vec}(\boldsymbol{\alpha}_{g^{(m)}}))^\top : m \neq l)^\top$, and $g_{\mathbf{i}_{-l}}^{-(l)} = (g_{i_m}^{(m)} : m \neq l)$. Using these notations, we define

$$\begin{aligned} Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \mathcal{G}_{-l}) &= \sum_{m \neq l} \sum_{i_m=1}^{N_m} \sum_{t=1}^T \left\{ Y_{i_1 \dots i_q, t} - \sum_{l=1}^q \lambda_{g_{i_l}}^{(l)} \sum_{k=1}^{N_l} w_{i_l k}^{(l)} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)} \right. \\ &\quad \left. - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} Y_{i_1 \dots i_q, (t-1)} - \sum_{l=1}^q \mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}}^{(l)} \right\}^2 \end{aligned} \quad (\text{A.1})$$

and $Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \mathcal{G}_{-l}) = E\{Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \mathcal{G}_{-l})\}$. By the definition, it readily follows that $Q(\boldsymbol{\xi}, \mathcal{G}) = \sum_{i_l=1}^{N_l} Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \mathcal{G}_{-l})$, where the $Q(\boldsymbol{\xi}, \mathcal{G})$ is the objective function defined in (9). For the tensor $\boldsymbol{\alpha} \in \mathbb{R}^{G_1 \times \dots \times G_q}$, we use $\boldsymbol{\alpha}_{g_{i_l}}^{(l)} \in \mathbb{R}^{G_1 \times \dots \times G_{l-1} \times G_{l+1} \times \dots \times G_q}$ to denote its subset by specifying the l th dimension being equal to $g_{i_l}^{(l)}$. Similarly, for the parameter tensor $\Theta \in \mathbb{R}^{N_1 \times \dots \times N_q \times m}$, when specifying the l th dimension as i_l , we obtain the subset tensor $\Theta_{i_l} \in \mathbb{R}^{N_1 \times \dots \times N_{l-1} \times N_{l+1} \times \dots \times N_q \times m}$. We define the pseudo distance

$$\begin{aligned} d_{i_l}(\widehat{\Theta}_{i_l}, \Theta_{i_l}) &= \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m=1}^{N_m} \left\| \widehat{\Theta}_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q} \right\|^2 \\ &= \sum_{m \neq l} \frac{1}{N_m} \sum_{i_m} \left\| \widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_m}^{(m)}}^{(m)} - \boldsymbol{\theta}_{g_{i_m}^{(m)}}^{(m)} \right\|^2 + \left\| \widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)}}^{(l)} \right\|^2 \\ &\quad + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m=1}^{N_m} \left| \widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} \right|^2. \end{aligned} \quad (\text{A.2})$$

For the parameter $\Theta_{i_1 \dots i_q}$, define the element-wise pseudo distance as

$$d_{i_1 \dots i_q}(\widehat{\Theta}_{i_1 \dots i_q}, \Theta_{i_1 \dots i_q}) = \sum_{l=1}^q \|\widehat{\boldsymbol{\theta}}_{g_{i_l}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}}^{(l)}\|^2 + |\widehat{\alpha}_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}}|^2. \quad (\text{A.3})$$

Let $\mathbb{Y}_t = \text{vec}(\mathcal{Y}_t) \in \mathbb{R}^{\prod_l N_l}$ and $\mathbb{E}_t = \text{vec}(\mathcal{E}_t) \in \mathbb{R}^{\prod_l N_l}$. Then (3) can be rewritten as $\mathbb{Y}_t = \mathbf{B}_0 \mathbb{Y}_{t-1} + \mathbf{c}_t + \mathbb{E}_t$, where

$$\mathbf{B}_0 = \sum_l \mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes (\mathbf{L}_{l,0} \mathbf{W}^{(l)}) \otimes \mathbf{I}_{N_{l+1}} \otimes \dots \otimes \mathbf{I}_{N_q} + \text{diag}\{\text{vec}(\mathcal{A}_0)\} \in \mathbb{R}^{\prod_l N_l \times \prod_l N_l}, \quad (\text{A.4})$$

$$\mathbf{c}_t = \sum_{l=1}^q \text{vec}(\mathbf{1}_{N_1} \circ \dots \circ \mathbf{1}_{N_{l-1}} \circ \boldsymbol{\beta}_{X_{i,t}}^{(l)0} \circ \mathbf{1}_{N_{l+1}} \circ \dots \circ \mathbf{1}_{N_q}) \stackrel{\text{def}}{=} \sum_l \mathbf{c}_t^{(l)} \in \mathbb{R}^{\prod_l N_l}, \quad (\text{A.5})$$

and $\mathbf{L}_{l,0}$, \mathcal{A}_0 , $\boldsymbol{\beta}_{X_{i,t}}^{(l)0}$ are the true values of the corresponding terms for $l \in [q]$, where $\boldsymbol{\beta}_{X_{i,t}}^{(l)} = (\mathbf{x}_{i,t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}}^{(l)} : 1 \leq i_l \leq N_l)^\top \in \mathbb{R}^{N_l}$ is defined in the main text. Recall that $E(\mathbf{c}_t) = \mathbf{0}$ because $E(\mathbf{x}_{i,t}^{(l)}) = \mathbf{0}$ for all $i_l \in [N_l]$ and $l \in [q]$, and by Assumption 5 that $\mathbb{Y}_0 = \mathbf{0}$, we have

$$\mathbb{Y}_t = \sum_{k=0}^t \mathbf{B}_0^k \mathbf{c}_{t-k} + \sum_{k=0}^t \mathbf{B}_0^k \mathbb{E}_{t-k} \stackrel{\text{def}}{=} \mathbb{Y}_t^c + \mathbb{Y}_t^e. \quad (\text{A.6})$$

Let $\boldsymbol{\Gamma} = \text{cov}(\mathbb{Y}_t)$. Denote $\tau_{\max} = \max_{i_1, \dots, i_q} \lambda_{\max}(\boldsymbol{\Sigma}_{i_1 \dots i_q})$ and $\boldsymbol{\Sigma}_{i_1 \dots i_q}$ is defined in Assumption 2.

GENERAL NOTATIONS. For a matrix $\mathbf{M} = (m_{ij}) \in \mathbb{R}^{m \times n}$, let $|\mathbf{M}|_e = (|m_{ij}|)_{i \in [m], j \in [n]}$. For two matrices $\mathbf{M}_1 = (m_{1ij}) \in \mathbb{R}^{m \times n}$ and $\mathbf{M}_2 = (m_{2ij}) \in \mathbb{R}^{m \times n}$, $\mathbf{M}_1 \preceq \mathbf{M}_2$ means that $m_{1ij} \leq m_{2ij}$ for all $i \in [m]$ and $j \in [n]$. In addition, define $\|\mathbf{M}\|$ as the largest singular value of \mathbf{M} , $\|\mathbf{M}\|_F = \text{tr}(\mathbf{M}\mathbf{M}^\top)^{1/2}$ as the Frobenius norm, $\|\mathbf{M}\|_\infty = \max_i \sum_j |m_{ij}|$, and $\|\mathbf{M}\|_{\max}$ as its largest absolute value. For a vector \mathbf{v} , denote the norm $\|\mathbf{v}\|_1 = \sum_i |v_i|$. Denote $\mathbf{e}_i^{(n)}$ as a vector of length n , whose i th element equals to 1 while others are 0.

Appendix B. Additional Discussion for the Case when $q = 2$

B.1 Estimation Details

As shown in Section 3.1, the estimator when $q = 2$ could be written as

$$\mathbf{M} = \begin{pmatrix} \mathbf{M}^{(1)} & \mathbf{M}^{(12)} & \mathbf{M}^{(1\alpha)} \\ \mathbf{M}^{(12)\top} & \mathbf{M}^{(2)} & \mathbf{M}^{(2\alpha)} \\ \mathbf{M}^{(1\alpha)\top} & \mathbf{M}^{(2\alpha)\top} & \mathbf{M}^\alpha \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \mathbf{b}^\alpha \end{pmatrix}, \quad (\text{A.7})$$

where the expressions are given as

$$\mathbf{M}_{g^{(1)}}^{(1)} = \sum_{t, g^{(2)}} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(1)\top} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(1)}, \quad \mathbf{M}^{(1)} = \text{diag}\left\{\mathbf{M}_{g^{(1)}}^{(1)} : g^{(1)} \in [G_1]\right\} \in \mathbb{R}^{G_1(p_1+1) \times G_1(p_1+1)},$$

$$\begin{aligned}
 \mathbf{M}_{g^{(1)}g^{(2)}}^{(12)} &= \sum_t \mathbb{X}_{g^{(1)}g^{(2)}t}^{(1)\top} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(2)}, \\
 \mathbf{M}^{(12)} &= \left(\mathbf{M}_{g^{(1)}g^{(2)}}^{(12)} : g^{(1)} \in [G_1], g^{(2)} \in [G_2] \right) \in \mathbb{R}^{G_1(p_1+1) \times G_2(p_2+1)}, \\
 \mathbf{M}_{g^{(1)}\mathcal{I}_{g^{(1)'g^{(2)}}}}^{(1\alpha)} &= \sum_t \mathbb{X}_{g^{(1)}g^{(2)}t}^{(1)\top} \mathbb{Y}_{g^{(1)}g^{(2)},(t-1)} I \left(g^{(1)} = g^{(1)'} \right), \quad \mathcal{I}_{g^{(1)'g^{(2)}}} = (g^{(2)} - 1)G_1 + g^{(1)'}, \\
 \mathbf{M}^{(1\alpha)} &= \left(\mathbf{M}_{g^{(1)}\mathcal{I}_{g^{(1)'g^{(2)}}}}^{(1\alpha)} : g^{(1)} \in [G_1], \mathcal{I}_{g^{(1)'g^{(2)}}} \in [G_1G_2] \right) \in \mathbb{R}^{G_1(p_1+1) \times G_1G_2}, \\
 \mathbf{M}_{g^{(2)}}^{(2)} &= \sum_{t,g^{(1)}} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(2)\top} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(2)}, \quad \mathbf{M}^{(2)} = \text{diag} \left\{ \mathbf{M}_{g^{(2)}}^{(2)} : g^{(2)} \in [G_2] \right\} \in \mathbb{R}^{G_2(p_2+1) \times G_2(p_2+1)}, \\
 \mathbf{M}_{g^{(2)}\mathcal{I}_{g^{(1)g^{(2)'}}} }^{(2\alpha)} &= \sum_t \mathbb{X}_{g^{(1)}g^{(2)}t}^{(2)\top} \mathbb{Y}_{g^{(1)g^{(2)'},(t-1)}} I \left(g^{(2)} = g^{(2)'} \right), \\
 \mathbf{M}^{(2\alpha)} &= \left(\mathbf{M}_{g^{(1)}\mathcal{I}_{g^{(1)g^{(2)'}}} }^{(2\alpha)} : g^{(2)} \in [G_2], \mathcal{I}_{g^{(1)g^{(2)'}}} \in [G_1G_2] \right) \in \mathbb{R}^{G_2(p_2+1) \times G_1G_2}, \\
 \mathbf{M}_{\mathcal{I}_{g^{(1)g^{(2)}}}\mathcal{I}_{g^{(1)'g^{(2)'}}} }^{\alpha} &= \sum_t \left\| \mathbb{Y}_{g^{(1)g^{(2)'},(t-1)}} \right\|^2 I \left(g^{(1)} = g^{(1)'}, g^{(2)} = g^{(2)'} \right), \\
 \mathbf{M}^{\alpha} &= \left(\mathbf{M}_{\mathcal{I}_{g^{(1)g^{(2)}}}\mathcal{I}_{g^{(1)'g^{(2)'}}} }^{\alpha} : \mathcal{I}_{g^{(1)g^{(2)}}} \in [G_1G_2], \mathcal{I}_{g^{(1)'g^{(2)'}}} \in [G_1G_2] \right) \in \mathbb{R}^{G_1G_2 \times G_1G_2}, \\
 \mathbf{b}_{g^{(1)}}^{(1)} &= \sum_{t,g^{(2)}} \mathbb{X}_{g^{(1)}g^{(2)}t}^{(1)\top} \mathbb{Y}_{g^{(1)g^{(2)},t}}, \quad \mathbf{b}^{(1)} = \left(\mathbf{b}_{g^{(1)}}^{(1)\top} : g^{(1)} \in [G_1] \right)^{\top} \in \mathbb{R}^{G_1(p_1+1)}, \\
 \mathbf{b}_g^{(2)} &= \sum_{t,g^{(1)}} \mathbb{X}_{g^{(1)g^{(2)}t} }^{(2)\top} \mathbb{Y}_{g^{(1)g^{(2)},t}}, \quad \mathbf{b}^{(2)} = \left(\mathbf{b}_{g^{(2)}}^{(2)\top} : g^{(2)} \in [G_2] \right)^{\top} \in \mathbb{R}^{G_2(p_2+1)}, \\
 \mathbf{b}_{\mathcal{I}_{g^{(1)g^{(2)}}}}^{\alpha} &= \sum_t \mathbb{Y}_{g^{(1)g^{(2)'},(t-1)}}^{\top} \mathbb{Y}_{g^{(1)g^{(2)},t}}, \quad \mathbf{b}^{\alpha} = \left(\mathbf{b}_{\mathcal{I}_{g^{(1)g^{(2)}}}}^{\alpha} : \mathcal{I}_{g^{(1)g^{(2)}}} \in [G_1G_2] \right)^{\top} \in \mathbb{R}^{G_1G_2}.
 \end{aligned}$$

In the above expressions,

$$\mathbb{X}_{g^{(1)}g^{(2)},t}^{(1)} = \left(\text{vec}(\mathbf{W}_1^{(\mathcal{R}_{g^{(1)'}, \cdot}^{(1)}, \cdot)} \mathbf{Y}_{t-1}^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})}), \mathbf{1}_{N_{2g^{(2)}}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_{g^{(1)'}, \cdot}^{(1)})} \right) \in \mathbb{R}^{(N_{1g^{(1)}}N_{2g^{(2)}}) \times (p_1+1)}, \quad (\text{A.8})$$

$$\mathbb{X}_{g^{(1)g^{(2)},t} }^{(2)} = \left(\text{vec}(\mathbf{Y}_{t-1}^{(\mathcal{R}_{g^{(1)'}, \cdot}^{(1)}, \cdot)} \mathbf{W}_2^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})}), (\mathbf{X}_t^{(2)})^{(\mathcal{R}_{g^{(2)'}, \cdot}^{(2)})} \otimes \mathbf{1}_{N_{1g^{(1)}}} \right) \in \mathbb{R}^{(N_{1g^{(1)}}N_{2g^{(2)}}) \times (p_2+1)}, \quad (\text{A.9})$$

and $\mathbf{X}_t^{(l)} = (\mathbf{x}_{1t}^{(l)}, \dots, \mathbf{x}_{N_1t}^{(l)})^{\top} \in \mathbb{R}^{N_l \times p_l}$, $N_{lg^{(l)}}^{(l)} = |\mathcal{R}_{g^{(l)}}^{(l)}|$ for $l = 1, 2$, and $\mathbb{Y}_{g^{(1)g^{(2)},t} } = \text{vec}(\mathcal{Y}_t^{(\mathcal{R}_{g^{(1)'}, \cdot}^{(1)}, \mathcal{R}_{g^{(2)}}^{(2)})}) \in \mathbb{R}^{|\mathcal{R}_{g^{(1)'}, \cdot}^{(1)}| |\mathcal{R}_{g^{(2)}}^{(2)}|}$. We then provide the iterative update algorithm when $q = 2$ in the following Algorithm A.1.

We remark here that the technical conditions in Section 4 can be simplified in this matrix case for easy understanding. For example, the covariates Assumption 4 can be rewritten as follow.

Assumption A.8. (DISTRIBUTION OF COVARIATES OF MATRIX CASE) *Recall that $\mathbf{x}_{it}^{(l)} \in \mathbb{R}^{p_l}$ is the covariate vector of the i th subject in the l th layer at time t . Assume $E(\mathbf{x}_{it}^{(1)}) = \mathbf{0}$*

Algorithm A.1 Estimation of the GTNAR Model When $q = 2$

- 1: **Input:** $\{\mathcal{Y}_t, \mathbf{X}_t^{(l)}, \mathbf{W}^{(l)}, G_l\}$ for $l = 1, 2$.
- 2: Obtain initial group memberships $\mathcal{G}_l^{[0]}$ according to Algorithm A.3 in Appendix F. Let $\{\boldsymbol{\xi}^{[k]}, \mathcal{G}_l^{[k]}\}$ be the estimators and memberships in the k th iteration.
- 3: Repeat STEP 1 and STEP 2 for $k = 1, 2, \dots$ until convergence.
 - STEP 1. Given $\{\mathcal{G}_l^{[k-1]}\}$ for all layers $l = 1, 2$, calculate $\boldsymbol{\xi}^{[k-1]} = (\mathbf{M}^{[k-1]})^{-1} \mathbf{b}^{[k-1]}$, where $\mathbf{M}^{[k-1]}$ and $\mathbf{b}^{[k-1]}$ are obtained from (13) with $\mathcal{G}_l^{[k-1]}$ s specified.
 - STEP 2. Given $\boldsymbol{\xi}^{[k-1]}$, update the memberships in layer $l = 1, 2$ sequentially by,

$$\begin{aligned}
 g_i^{(1)[k]} = & \arg \min_{g_i^{(1)} \in [G_1]} \sum_{j=1}^{N_2} \sum_{t=1}^T \left\{ Y_{ij,t} \right. \\
 & - \left(\lambda_{g_i^{(1)}}^{(1)[k]} \sum_{m=1}^{N_1} \frac{a_{im}^{(1)}}{n_{1i}} Y_{mj,(t-1)} + \lambda_{g_j^{(2)[k-1]}}^{(2)[k]} \sum_{m=1}^{N_2} \frac{a_{jm}^{(2)}}{n_{2j}} Y_{im,(t-1)} \right) \\
 & \left. - \alpha_{g_i^{(1)} g_j^{(2)[k-1]}^{[k]} Y_{ij,(t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)[k]} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)[k-1]}}^{(2)[k]} \right) \right\}^2 \tag{A.10}
 \end{aligned}$$

and

$$\begin{aligned}
 g_j^{(2)[k]} = & \arg \min_{g_j^{(2)} \in [G_2]} \sum_{i=1}^{N_1} \sum_{t=1}^T \left\{ Y_{ij,t} \right. \\
 & - \left(\lambda_{g_i^{(1)[k]}}^{(1)[k]} \sum_{m=1}^{N_1} \frac{a_{im}^{(1)}}{n_{1i}} Y_{mj,(t-1)} + \lambda_{g_j^{(2)}}^{(2)[k]} \sum_{m=1}^{N_2} \frac{a_{jm}^{(2)}}{n_{2j}} Y_{im,(t-1)} \right) \\
 & \left. - \alpha_{g_i^{(1)[k]} g_j^{(2)}^{[k]} Y_{ij,(t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)[k]}}^{(1)[k]} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)[k]} \right) \right\}^2 \tag{A.11}
 \end{aligned}$$

- 4: **Output:** Final estimator and memberships: $\hat{\boldsymbol{\xi}} = \boldsymbol{\xi}^{[K]}$ and $\hat{\mathcal{G}} = \{\mathcal{G}_l^{[K]} : l = 1, 2\}$. Here K is the final number of iteration rounds.
-

for all $i \in [N_1]$ and $t \in [T]$. Let $\boldsymbol{\eta}_1 \in \mathbb{R}^{p_1}$ be a constant vector satisfying $\|\boldsymbol{\eta}_1\| \leq c$, where c is a positive constant. Define $\mathbf{x}_t^{(1)\eta} = (\mathbf{x}_{it}^{(1)\top} \boldsymbol{\eta}_1 : i \in [N_1])^\top \in \mathbb{R}^{N_1}$. Assume $\mathbf{x}^{(1)\eta} = (\mathbf{x}_t^{(1)\eta\top} : 0 \leq t \leq T)^\top \in \mathbb{R}^{N_1(T+1)}$ satisfies the K -convex concentration property for some constant K according to Definition 1.

B.2 Model Inference

In this subsection, we provide an estimator to the asymptotic covariance when $q = 2$. This procedure can be easily extended to $q > 2$. With the parameter estimator $\widehat{\boldsymbol{\Theta}} = (\widehat{\boldsymbol{\Theta}}_{ij} = (\widehat{\boldsymbol{\theta}}_{g_i}^{(1)\top}, \widehat{\boldsymbol{\theta}}_{g_j}^{(2)\top}, \widehat{\alpha}_{g_i(1)g_j(2)})^\top)$, we first estimate σ^2 as follows

$$\widehat{\sigma}^2 = \frac{1}{N_1 N_2 T} \sum_{i=1} \sum_{j=1} \sum_{t=1} (Y_{ij,t} - \mathcal{X}_{ij,t}^\top \widehat{\boldsymbol{\Theta}}_{ij})^2, \quad (\text{A.12})$$

where $\mathcal{X}_{ij,t}$ is defined in Assumption 2. Next, we estimate \mathbf{M}^0 by $\widehat{\mathbf{M}}$, where $\widehat{\mathbf{M}}$ is obtained by plugging estimated parameters $\widehat{\boldsymbol{\Theta}}$ into the expression in (13). In the following theorem, we show that the covariance estimator is consistent.

Theorem 10. *Suppose Assumption 1–5 and assume that there exists n , such that $c_1 n \leq \min_l N_l \leq \max_l N_l \leq c_2 n$ for some constants $c_1, c_2 > 0$. When $G_1 = G_{1,0}$ $G_2 = G_{2,0}$, assume $\{\log(N_1 N_2)\}^2 / T \rightarrow 0$, then the following holds,*

$$\widehat{\sigma}^2 \rightarrow_p \sigma^2, \quad (n^2 T)^{-1} \widehat{\sigma}^2 \boldsymbol{\eta}^\top (\widehat{\mathbf{M}})^{-1} \boldsymbol{\eta} \rightarrow_p (n^2 T)^{-1} \sigma^2 \boldsymbol{\eta}^\top (\mathbf{M}^0)^{-1} \boldsymbol{\eta},$$

where $\boldsymbol{\eta}$ is defined in Theorem 7.

The proof of Theorem 10 can be found in Appendix B.2.1. Theorem 10 indicates that we can obtain a consistent estimator for the asymptotic variance by plugging in the estimators $\widehat{\boldsymbol{\Theta}}$ and the consistent estimator $\widehat{\sigma}^2$. This assures a valid statistical inference procedure. We next present a number of simulation studies to examine the finite sample performances of the model estimation and inference procedures.

B.2.1 PROOF OF THEOREM 10

To prove Theorem 10, we prove the following two statements separately,

$$\widehat{\sigma}^2 \rightarrow_p \sigma^2, \quad (\text{A.13})$$

$$(n^2 T)^{-1} \widehat{\mathbf{M}}_{nT} \rightarrow_p (n^2 T)^{-1} \mathbf{M}^0. \quad (\text{A.14})$$

(1) Proof of (A.13)

Note that $\sigma^2 = (N_1 N_2 T)^{-1} \sum_{i,j,t} E(\varepsilon_{ij,t}^2) = (N_1 N_2 T)^{-1} E\{\sum_{i,j} Q_{ij}(\boldsymbol{\Theta}_{ij}^0)\} = (N_1 N_2 T)^{-1} Q^*(\boldsymbol{\Theta}^0)$, and $\widehat{\sigma}^2 = (N_1 N_2 T)^{-1} \sum_{i,j} Q_{ij}(\widehat{\boldsymbol{\Theta}}_{ij}) = (N_1 N_2 T)^{-1} Q(\widehat{\boldsymbol{\Theta}})$. To prove (A.13), it suffices to show that $(N_1 N_2 T)^{-1} |Q(\widehat{\boldsymbol{\Theta}}) - Q^*(\boldsymbol{\Theta}^0)| = o_p(1)$. Further note that

$$\frac{1}{N_1 N_2 T} |Q(\widehat{\boldsymbol{\Theta}}) - Q^*(\boldsymbol{\Theta}^0)| \leq \frac{1}{N_1 N_2 T} \left\{ |Q(\widehat{\boldsymbol{\Theta}}) - Q^*(\widehat{\boldsymbol{\Theta}})| + |Q^*(\widehat{\boldsymbol{\Theta}}) - Q^*(\boldsymbol{\Theta}^0)| \right\},$$

hence we bound the two terms on the right hand side.

(1.1) Order of $(N_1 N_2 T)^{-1} |Q(\widehat{\Theta}) - Q^*(\widehat{\Theta})|$.

By Lemma 17, we have that

$$\sup_{\|\Theta\|_{\max} \leq R} \left| \frac{1}{N_1 N_2 T} \{Q(\Theta) - Q^*(\Theta)\} \right| = O_p \left\{ T^{-1/2} (m + \log(N_1 N_2)) \right\},$$

where $m = p_1 + p_2 + 3$. Thus when $\{m + \log(N_1 N_2)\}/\sqrt{T} \rightarrow 0$, we have $(N_1 N_2 T)^{-1} |Q(\widehat{\Theta}) - Q^*(\widehat{\Theta})| = o_p(1)$.

(1.2) Order of $(N_1 N_2 T)^{-1} |Q^*(\widehat{\Theta}) - Q^*(\Theta^0)|$

By Lemma 23 we have $(N_1 N_2 T)^{-1} \{Q^*(\widehat{\Theta}) - Q^*(\Theta^0)\} \leq \tau_{\max} d(\widehat{\Theta}, \Theta^0)$. Since τ_{\max} is bounded due to Lemma 33 and $d(\widehat{\Theta}, \Theta^0) = O_p\{T^{-1/2}(m + \log(N_1 N_2))\}$ according to Theorem 2, we have that $(N_1 N_2 T)^{-1} |Q^*(\widehat{\Theta}) - Q^*(\Theta^0)| = O_p\{T^{-1}(\log(N_1 N_2)^2)\} = o_p(1)$ when $(\log N_1 N_2)^2/T \rightarrow 0$. This condition can be implied by $\{m + \log(N_1 N_2)\}/\sqrt{T} \rightarrow 0$.

From both (1.1) and (1.2), we have that $(N_1 N_2 T)^{-1} |Q(\widehat{\Theta}) - Q^*(\Theta^0)| = o_p(1)$ when $\{m + \log(N_1 N_2)\}/\sqrt{T} \rightarrow 0$. This implies that $\widehat{\sigma}^2 \rightarrow_p \sigma^2$ under the same condition. This finishes the proof of (A.13).

(2) Proof of (A.14)

Recall the expression of \mathbf{M} defined in (A.7) when $q = 2$, we need to show that the nine elements in $(n^2 T)^{-1} \widehat{\mathbf{M}}_{nT}$ are consistent in an analogous manner. Hence, we take the first element, i.e. $(n^2 T)^{-1} \widehat{\mathbf{M}}^{(1)} \in \mathbb{R}^{G_1(p_1+1) \times G_1(p_1+1)}$, as an example to show the consistency, where n is defined in Theorem 10. Other terms could be proved similarly.

Recall that $\mathbf{M}^{(1)} = \text{diag}\{\mathbf{M}_{g^{(1)}}^{(1)} : g^{(1)} \in [G_1]\}$ and $\mathbf{M}_{g^{(1)}}^{(1)} = \sum_{t, g^{(2)}} \mathbb{X}_{g^{(1)}g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}g^{(2)}, t}$, we show that for any $g^{(1)} \in [G_{1,0}]$, it holds that $|(n^2 T)^{-1} \widehat{\mathbf{M}}_{g^{(1)}}^{(1)} - \lim_{n, T \rightarrow \infty} (n^2 T)^{-1} E(\widehat{\mathbf{M}}_{g^{(1)}}^{(1)})| = o_p(1)$. Recall that

$$\begin{aligned} \mathbb{X}_{g^{(1)}g^{(2)}, t} &= \left(\text{vec}(\mathbf{W}_1^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot, \mathcal{R}_{g^{(2)}}^{(2)})}, \mathbf{Y}_{t-1}^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})}, \mathbf{1}_{N_{2g^{(2)}}} \otimes \mathbf{X}_t^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})} \right) \\ &= \left((\mathbf{I}_{N_{2g^{(2)}}} \otimes \mathbf{W}_1^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot, \mathcal{R}_{g^{(2)}}^{(2)})}) \mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)})}, \mathbf{1}_{N_{2g^{(2)}}} \otimes \mathbf{X}_t^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})} \right) \\ &\stackrel{\text{def}}{=} \left(\mathbb{W}_{g^{(1)}g^{(2)}}^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot, \mathcal{R}_{g^{(2)}}^{(2)})}, \mathbf{1}_{N_{2g^{(2)}}} \otimes \mathbf{X}_t^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})} \right) \in \mathbb{R}^{(N_{1g^{(1)}} N_{2g^{(2)}}) \times (p_1+1)}, \end{aligned}$$

where $\mathbb{W}_{g^{(1)}g^{(2)}} \stackrel{\text{def}}{=} \mathbf{I}_{N_{2g^{(2)}}} \otimes \mathbf{W}_1^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot, \mathcal{R}_{g^{(2)}}^{(2)})} \in \mathbb{R}^{N_{2g^{(2)}} N_{1g^{(1)}} \times N_{2g^{(2)}} N_1}$ and $\mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)})} \stackrel{\text{def}}{=} \text{vec}(\mathbf{Y}_{t-1}^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})}) \in \mathbb{R}^{N_{2g^{(2)}} N_1}$. Hence, it is equivalent to prove that

$$\begin{aligned} &(n^2 T)^{-1} \sum_{t, g^{(2)}} \mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)}) \top} \mathbb{W}_{g^{(1)}g^{(2)}}^\top \mathbb{W}_{g^{(1)}g^{(2)}} \mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)})} \\ &\rightarrow_p \lim_{n, T \rightarrow \infty} (n^2 T)^{-1} \sum_{t, g^{(2)}} E(\mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)}) \top} \mathbb{W}_{g^{(1)}g^{(2)}}^\top \mathbb{W}_{g^{(1)}g^{(2)}} \mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)})}), \quad (\text{A.15}) \\ &(n^2 T)^{-1} \sum_{t, g^{(2)}} \mathbb{Y}_{t-1}^{(\mathcal{R}_{g^{(2)}}^{(2)}) \top} \mathbb{W}_{g^{(1)}g^{(2)}}^\top \mathbf{1}_{N_{2g^{(2)}}} \otimes \mathbf{X}_t^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})} \end{aligned}$$

$$\rightarrow_p \lim_{n, T \rightarrow \infty} (n^2 T)^{-1} \sum_{t, g^{(2)}} E(\mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbf{1}_{N_{2g^{(2)}}} \otimes \mathbf{X}_t^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})}). \quad (\text{A.16})$$

We prove (A.15) for example, and (A.16) could be proved similarly. Similar as the decomposition in (A.6), write $\mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)}} = \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, c}^{(2)}} + \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, e}^{(2)}}$, where $\mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, c}^{(2)}} = \sum_{k=0}^{t-1} (\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}})^k \mathbf{c}_{t-k}^{\mathcal{R}_{g^{(2)}}^{(2)}}$, $\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}} = \mathbf{I}_{N_{2g^{(2)}}} \otimes \mathbf{L}_0 \mathbf{W}_1 + \mathbf{G}_0^{(\mathcal{R}_{g^{(2)}}, \mathcal{R}_{g^{(2)}}^{(2)})\top} \mathbf{W}_2^{(\mathcal{R}_{g^{(2)}}, \mathcal{R}_{g^{(2)}}^{(2)})} \otimes \mathbf{I}_{N_1} \in \mathbb{R}^{N_1 N_{2g^{(2)}} \times N_{2g^{(2)}} N_1}$, and $\mathbf{c}_t^{\mathcal{R}_{g^{(2)}}^{(2)}} \in \mathbb{R}^{N_1 N_{2g^{(2)}}}$ is a similar modification of \mathbf{c}_t . $\mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, e}^{(2)}}$ could be represented analogously. Then the left hand side of (A.15) could be separated as

$$\begin{aligned} & \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &= \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, c}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, c}^{(2)}} + 2 \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, c}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, e}^{(2)}} \\ &+ \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, e}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}, e}^{(2)}} \end{aligned}$$

We take the first term as an example. Let $\mathbf{N}_{g^{(1)g^{(2)}}} = \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \in \mathbb{R}^{N_{2g^{(2)}} N_1 \times N_{2g^{(2)}} N_1}$, then we have

$$\begin{aligned} & \frac{1}{n^2 T} \sum_{t, g^{(2)}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &= \frac{1}{n^2 T} \sum_{t, g^{(2)}} \sum_{l_1=0}^{T, G_2} \sum_{l_2=0}^{t-1} \mathbf{c}_{l_1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} (\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}})^{t-l_1\top} \mathbf{N}_{g^{(1)g^{(2)}}} (\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}})^{t-l_2} \mathbf{c}_{l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &= \frac{1}{n} \sum_{g^{(2)}} \sum_{l_1=0}^{T-1} \sum_{l_2=0}^{T-1} \mathbf{c}_{l_1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \left\{ \frac{1}{nT} \sum_{t=\max(1, l_1, l_2)}^T (\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}})^{t-l_1\top} \mathbf{N}_{g^{(1)g^{(2)}}} (\mathbf{B}_0^{\mathcal{R}_{g^{(2)}}^{(2)}})^{t-l_2} \right\} \mathbf{c}_{l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{g^{(2)}} \sum_{l_1=0}^{T-1} \sum_{l_2=0}^{T-1} \mathbf{c}_{l_1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbf{E}_{l_1, l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{c}_{l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}}. \end{aligned}$$

Define $\mathbf{c}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} = (\mathbf{c}_0^{\mathcal{R}_{g^{(2)}}^{(2)\top}}, \dots, \mathbf{c}_{T-1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}})^\top$, and $\mathbf{E}_{l_1, l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}} = (\mathbf{E}_{l_1, l_2}^{\mathcal{R}_{g^{(2)}}^{(2)}})$, $0 \leq l_1, l_2 \leq T-1$, then we have

$$\begin{aligned} & \frac{1}{n^2 T} \sum_{t, g^{(2)}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}^{(2)}} = \frac{1}{n} \sum_{g^{(2)}} \mathbf{c}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)\top}} \mathbf{E}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{c}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &= \frac{1}{n} \sum_{g^{(2)}} \mathbf{c}_{g^{(2)}, (1)\top}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{E}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{c}_{g^{(2)}, (1)}^{\mathcal{R}_{g^{(2)}}^{(2)}} + \frac{2}{n} \sum_{g^{(2)}} \mathbf{c}_{g^{(2)}, (1)\top}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{E}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{c}_{g^{(2)}, (2)}^{\mathcal{R}_{g^{(2)}}^{(2)}} \\ &+ \frac{1}{n} \sum_{g^{(2)}} \mathbf{c}_{g^{(2)}, (2)\top}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{E}_{g^{(2)}}^{\mathcal{R}_{g^{(2)}}^{(2)}} \mathbf{c}_{g^{(2)}, (2)}^{\mathcal{R}_{g^{(2)}}^{(2)}}. \end{aligned}$$

Then we could follow the proof of Lemma 21 to show the result. Since

$$n^{-2} \sum_{g^{(2)}} \mathbf{1}_{N_{2g^{(2)}} N_1} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbf{1}_{N_{2g^{(2)}} N_1} = O(1)$$

under Assumption that there exists n that $c_1 n \leq \min_l N_l \leq \max_l N_l \leq c_2 n$, we would have

$$\begin{aligned} & P \left\{ \left| \frac{1}{n^2 T} \sum_{t, g^{(2)}} \left\{ \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}} - E(\mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}} \mathbb{W}_{g^{(1)g^{(2)}}}^\top \mathbb{W}_{g^{(1)g^{(2)}}} \mathbb{Y}_{t-1}^{\mathcal{R}_{g^{(2)}}}) \right\} \right| \geq u \right\} \\ & \leq C \exp \left\{ -c \min(Tu^2, \sqrt{T}u) \right\}. \end{aligned}$$

This would lead to the statement (A.15). Then, (A.16) can be shown similarly. Therefore, the consistency in (A.14) could be obtained when $T \rightarrow \infty$. Overall, the consistency of the asymptotic covariance estimation in Theorem 10 could be obtained by Step (1) and Step (2).

B.3 Existence of Asymptotic Covariance

We illustrate the asymptotic covariance matrix \mathbf{M}^0 in Theorem 7 exists when $q = 2$, which can be directly extended to the general q case. For notation simplicity, in this subsection, we use g and h to represent group in the first and second dimension, respectively.

Assumption A.9. Recall that $\mathbb{Y}_t = \text{vec}(\mathbf{Y}_t) \in \mathbb{R}^{N_1 N_2}$. Denote $\mathbb{Y}_{gh,t} = \text{vec}(\mathbf{Y}_t^{(\mathcal{R}_g^{(1)}, \mathcal{R}_h^{(2)})})$. For any $g_1, g_2 \in [G_1]$ and $h_1, h_2 \in [G_2]$, denote $\Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E} = E(\mathbb{Y}_{g_1 h_1, t}^{\mathcal{E}} \mathbb{Y}_{g_2 h_2, t}^{\mathcal{E}\top})$, $\Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X} = E(\mathbb{Y}_{g_1 h_1, t}^{\mathcal{E}} \mathbb{Y}_{g_2 h_2, t}^{\mathbf{X}\top})$. Denote $\mathcal{W}_{gh}^1 = (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g^{(1)}, \cdot)}) \in \mathbb{R}^{N_{1g} N_{2h} \times N_1 N_{2h}}$, and denote $\mathcal{W}_{gh}^2 = (\mathbf{I}_{N_{1g}} \otimes (\mathbf{W}_2^\top)^{(\cdot, \mathcal{R}_h^{(2)})}) \in \mathbb{R}^{N_{1g} N_{2h} \times N_1 N_{2h}}$. Assume the following limits exist,

$$\begin{aligned} \kappa_{g_1 h_1, g_2 h_2}^{11} &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E} \mathcal{W}_{g_1 h_2}^{1\top}), \quad \nu_{g_1 h_1, g_2 h_2}^{11} = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X} \mathcal{W}_{g_1 h_2}^{1\top}), \\ \kappa_{g_1 h_1, g_2 h_2}^{22} &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^2 \Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E} \mathcal{W}_{g_1 h_1}^{2\top}), \quad \nu_{g_1 h_1, g_2 h_2}^{22} = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^2 \Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X} \mathcal{W}_{g_1 h_1}^{2\top}), \\ \kappa_{g_1 h_1, g_2 h_2}^{12} &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E} \mathcal{W}_{g_1 h_1}^{2\top}), \quad \nu_{g_1 h_1, g_2 h_2}^{12} = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X} \mathcal{W}_{g_1 h_1}^{2\top}), \\ \kappa_{g_1 h_1, g_2 h_2}^1 &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E}), \quad \nu_{g_1 h_1, g_2 h_2}^1 = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^1 \Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X}), \\ \kappa_{g_1 h_1, g_2 h_2}^2 &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^2 \Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E}), \quad \nu_{g_1 h_1, g_2 h_2}^2 = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\mathcal{W}_{g_1 h_1}^2 \Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X}), \\ \kappa_{g_1 h_1, g_2 h_2}^\alpha &= \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\Sigma_{g_1 h_1, g_2 h_2}^\mathcal{E}), \quad \nu_{g_1 h_1, g_2 h_2}^\alpha = \lim_{n \rightarrow \infty} n^{-2} \text{tr}(\Sigma_{g_1 h_1, g_2 h_2}^\mathbf{X}), \end{aligned}$$

where $\Sigma_{h_1, h_2}^\mathcal{E}, \Sigma_{h_1, h_2}^\mathbf{X} \in \mathbb{R}^{N_{1N_2 h_1} \times N_{1N_2 h_2}}$, $\Sigma_{g_1, g_2}^\mathcal{E}, \Sigma_{g_1, g_2}^\mathbf{X} \in \mathbb{R}^{N_{1g_1} N_2 \times N_{1g_2} N_2}$, and $\Sigma_{g_1, g_2 h_2}^\mathcal{E}, \Sigma_{g_1, g_2 h_2}^\mathbf{X} \in \mathbb{R}^{N_{1g_1} N_2 \times N_{1g_2} N_{2h_2}}$, $\Sigma_{g_1 h_1, g_2, \cdot}^\mathcal{E}, \Sigma_{g_1 h_1, g_2, \cdot}^\mathbf{X} \in \mathbb{R}^{N_{1g_1} N_{2h_2} \times N_{1g_1} N_2}$. Further assume that $\Sigma_{X_1} = E(\mathbf{x}_{it}^{(1)} \mathbf{x}_{it}^{(1)\top}) \in \mathbb{R}^{p_1 \times p_1}$, $\Sigma_{X_2} = E(\mathbf{x}_{jt}^{(2)} \mathbf{x}_{jt}^{(2)\top}) \in \mathbb{R}^{p_2 \times p_2}$ exist.

Lemma 11. Under Assumptions A.9, the asymptotic covariance matrix $\mathbf{M}^0 = \lim_{n, T \rightarrow \infty} \mathbf{M}_{n, T}^0$ can be written as

$$\mathbf{M}^0 = \begin{pmatrix} \Xi_1 & \Xi_{12} & \Xi_{1\alpha} \\ \Xi_{12}^\top & \Xi_2 & \Xi_{2\alpha} \\ \Xi_{1\alpha}^\top & \Xi_{2\alpha}^\top & \Xi_\alpha \end{pmatrix}.$$

where

$$\begin{aligned}
 \Xi_1 &= \text{diag} \left(\left(\begin{array}{cc} s_{11,g} & \mathbf{0}^\top \\ \mathbf{0} & \Sigma_{X_1} \end{array} \right) : g \in [G_1] \right), \quad \Xi_{2,h} = \text{diag} \left(\left(\begin{array}{cc} s_{22,h} & \mathbf{0}^\top \\ \mathbf{0} & \Sigma_{X_2} \end{array} \right) : h \in [G_2] \right), \\
 \Xi_{12} &= \left(\left(\begin{array}{cc} s_{12,gh} & \mathbf{0}^\top \\ \mathbf{0} & \mathbf{0} \end{array} \right) : g \in [G_1], h \in [G_2] \right), \\
 \Xi_{1\alpha} &= \left((s_{31,gh} I(g = g'), \mathbf{0}^\top) : g \in [G_1], \mathcal{I}_{g'h} \in [G_1 G_2] \right), \\
 \Xi_{2\alpha} &= \left((s_{32,gh} I(h = h'), \mathbf{0}^\top) : h \in [G_2], \mathcal{I}_{gh'} \in [G_1 G_2] \right), \\
 \Xi_\alpha &= (s_{4,gh} I(g = g', h = h') : \mathcal{I}_{gh} \in [G_1 G_2], \mathcal{I}_{g'h'} \in [G_1 G_2]), \tag{A.17}
 \end{aligned}$$

In (A.17), $s_{11,g} = \sum_h \nu_{gh,gh}^{11} + \kappa_{gh,gh}^{11}$, $s_{22,h} = \sum_g \nu_{gh,gh}^{22} + \kappa_{gh,gh}^{22}$, $s_{12,gh} = \nu_{gh,gh}^{12} + \kappa_{gh,gh}^{12}$, $s_{31,gh} = \nu_{gh,gh}^1 + \kappa_{gh,gh}^1$, $s_{32,gh} = \nu_{gh,gh}^2 + \kappa_{gh,gh}^2$ and $s_{4,gh} = \nu_{gh,gh}^\alpha + \kappa_{gh,gh}^\alpha$, with all the constants defined in Assumption A.9.

Proof Recall the expression of \mathbf{M} when $q = 2$ is defined in (A.7),

$$\mathbf{M} = \begin{pmatrix} \mathbf{M}^{(1)} & \mathbf{M}^{(12)} & \mathbf{M}^{(1\alpha)} \\ \mathbf{M}^{(12)\top} & \mathbf{M}^{(2)} & \mathbf{M}^{(2\alpha)} \\ \mathbf{M}^{(1\alpha)\top} & \mathbf{M}^{(2\alpha)\top} & \mathbf{M}^\alpha \end{pmatrix}.$$

Due to the similar format of the terms, we derive the expression of

$$\begin{aligned}
 \mathbf{s}_1 &\stackrel{\text{def}}{=} \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} E(\mathbf{M}^{(1)}), \quad \mathbf{s}_2 \stackrel{\text{def}}{=} \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} E(\mathbf{M}^{(12)}), \\
 \mathbf{s}_3 &\stackrel{\text{def}}{=} \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} E(\mathbf{M}^{(1\alpha)}), \quad \mathbf{s}_4 \stackrel{\text{def}}{=} \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} E(\mathbf{M}^\alpha)
 \end{aligned}$$

in the following four steps.

(1). **Expression of \mathbf{s}_1 .**

We prove that $\mathbf{s}_1 = \text{diag}((s_{11,g}, \mathbf{s}_{12,g}^\top; \mathbf{s}_{12,g}, \mathbf{s}_{13}) : g \in [G_1])$ with $s_{11,g} = \sum_h s_{11,gh} = \sum_h \nu_{gh,gh}^{11} + \kappa_{gh,gh}^{11}$, $\mathbf{s}_{12,g} = \sum_h \mathbf{s}_{12,gh} = \mathbf{0}$ and $\mathbf{s}_{13,g} = \Sigma_{X_1}$, which is defined in Assumption A.9. Recall that $\mathbf{M}^{(1)} = \text{diag}(\mathbf{M}_g^{(1)} : g \in [G_1]) \in \mathbb{R}^{G_1(p_1+1) \times G_1(p_1+1)}$ with $\mathbf{M}_g^{(1)} = \sum_{t,h} \mathbb{X}_{gh,t}^{(1)\top} \mathbb{X}_{gh,t}^{(1)}$, where $\mathbb{X}_{gh,t}^{(1)}$ is defined in (A.8). By (A.8), denote $\mathbb{Y}_{\cdot,h,(t-1)} = \text{vec}(\mathbf{Y}_{t-1}^{(\cdot, \mathcal{R}_h^{(2)})})$, we have

$$\begin{aligned}
 s_{11,gh} &= \lim_{n \rightarrow \infty} (n^2 T)^{-1} \sum_t E \{ \mathbb{Y}_{\cdot,h,(t-1)}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g, \cdot)})^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g, \cdot)}) \mathbb{Y}_{\cdot,h,(t-1)} \}, \\
 \mathbf{s}_{12,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E \{ \mathbb{Y}_{\cdot,h,(t-1)}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g, \cdot)})^\top (\mathbf{1}_{N_{2h}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_g, \cdot)}) \}, \\
 \mathbf{s}_{13,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E \{ (\mathbf{1}_{N_{2h}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_g, \cdot)})^\top (\mathbf{1}_{N_{2h}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_g, \cdot)}) \}.
 \end{aligned}$$

Since $\mathbb{Y}_{\cdot,h,(t-1)}$ is independent with $\mathbf{X}_t^{(1)}$, we have $\mathbf{s}_{12} = \mathbf{0}$. Further note that

$$s_{11,gh} = \lim_{n \rightarrow \infty} (n^2 T)^{-1} \sum_t [\text{tr}(\mathcal{W}_{gh}^1 \Sigma_{\cdot,h,\cdot}^{\mathbf{X}} \mathcal{W}_{gh}^{1\top}) + \text{tr}(\mathcal{W}_{gh}^1 \Sigma_{\cdot,h,\cdot}^{\mathcal{E}} \mathcal{W}_{gh}^{1\top})]$$

$$= \lim_{n \rightarrow \infty} n^{-2} [\text{tr}(\mathcal{W}_{gh}^1 \boldsymbol{\Sigma}_{\cdot,h,h}^{\mathbf{X}} \mathcal{W}_{gh}^{1\top}) + \text{tr}(\mathcal{W}_{gh}^1 \boldsymbol{\Sigma}_{\cdot,h,h}^{\mathcal{E}} \mathcal{W}_{gh}^{1\top})] = \nu_{gh,gh}^{11} + \kappa_{gh,gh}^{11}.$$

Recall that $\mathbf{x}_{it}^{(1)}$ s are independent and identically distributed covariates across all $i \in [N_1]$ and $t \in [T]$, denote the covariance as $\boldsymbol{\Sigma}_{X_1} = E(\mathbf{x}_{it}^{(1)} \mathbf{x}_{it}^{(1)\top}) \in \mathbb{R}^{p_1 \times p_1}$. Therefore we have $\mathbf{s}_{13,gh} = \boldsymbol{\Sigma}_{X_1}$.

(2). Expression of \mathbf{s}_2 .

Next, we prove that $\mathbf{s}_2 = ((s_{21,gh}, \mathbf{s}_{22,gh}; \mathbf{s}_{22,gh}^\top, \mathbf{s}_{23,gh}) : g \in [G_1], h \in [G_2])$ with $\mathbf{s}_{21,gh} = \nu_{gh,gh}^{12} + \kappa_{gh,gh}^{12}$, $\mathbf{s}_{22,gh} = \mathbf{0}$ and $\mathbf{s}_{23,gh} = \mathbf{0}$. Since $\mathbf{M}^{(12)} = (\mathbf{M}_{gh}^{(12)} : g \in [G_1], h \in [G_2])$, we derive the expression of $\lim_{n,T \rightarrow \infty} (n^2 T)^{-1} E(\mathbf{M}_{gh}^{(12)}) = \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E(\mathbb{X}_{gh,t}^{(1)\top} \mathbb{X}_{gh,t}^{(2)})$. By the definition of $\mathbb{X}_{gh,t}^{(1)}$ and $\mathbb{X}_{gh,t}^{(2)}$ in (A.8) and (A.9), we have

$$\begin{aligned} s_{21,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{\mathbb{Y}_{\cdot,h,(t-1)}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g^{(1)}, \cdot)})^\top (\mathbf{I}_{N_{1g}} \otimes \mathbf{W}_2^{(\cdot, \mathcal{R}_h^{(2)})^\top}) \mathbb{Y}_{g \cdot, (t-1)}\}, \\ \mathbf{s}_{22,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{\mathbb{Y}_{\cdot,h,(t-1)}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g^{(1)}, \cdot)})^\top ((\mathbf{X}_t^{(2)})^{(\mathcal{R}_h^{(2)}, \cdot)} \otimes \mathbf{1}_{N_{1g}})\}, \\ \mathbf{s}_{23,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{((\mathbf{X}_t^{(2)})^{(\mathcal{R}_h^{(2)}, \cdot)} \otimes \mathbf{1}_{N_{1g}})^\top (\mathbf{1}_{N_{2h}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_g^{(1)}, \cdot)})\}. \end{aligned}$$

Since $\mathbb{Y}_{\cdot,h,(t-1)}$ is independent with $\mathbf{X}_t^{(2)}$, we have $\mathbf{s}_{22} = \mathbf{0}$, and due to the independence between $\mathbf{X}_t^{(1)}$ and $\mathbf{X}_t^{(2)}$, we have $\mathbf{s}_{23} = \mathbf{0}$. Note that

$$s_{21,gh} = \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t [\text{tr}(\mathcal{W}_{gh}^1 \boldsymbol{\Sigma}_{\cdot,h,g}^{\mathbf{X}} \mathcal{W}_{gh}^{2\top}) + \text{tr}(\mathcal{W}_{gh}^1 \boldsymbol{\Sigma}_{\cdot,h,g}^{\mathcal{E}} \mathcal{W}_{gh}^{2\top})] = \nu_{gh,gh}^{12} + \kappa_{gh,gh}^{12},$$

we obtain the final expression.

(3). Expression of \mathbf{s}_3 .

We show that $\mathbf{s}_3 = ((s_{31,gh} I(g = g'), \mathbf{s}_{32,gh}) : g \in [G_1], \mathcal{I}_{g'h} \in [G_1 G_2])$ with $\mathcal{I}_{g'h} = (h-1)G_1 + g'$, where $s_{31,gh} = \nu_{gh,gh}^1 + \kappa_{gh,gh}^1$ and $\mathbf{s}_{32,gh} = \mathbf{0}$. Since $\mathbf{M}^{(1\alpha)} = (\mathbf{M}_{g\mathcal{I}_{g'h}}^{(1\alpha)} : g \in [G_1], \mathcal{I}_{g'h} \in [G_1 G_2])$ and $\mathbf{M}_{g\mathcal{I}_{g'h}}^{(1\alpha)} = \sum_t \mathbb{X}_{gh,t}^{(1)\top} \mathbb{Y}_{gh,(t-1)} I(g = g')$ with $\mathcal{I}_{g'h} = (h-1)G_1 + g'$, we derive the expression of $\lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t \mathbb{X}_{gh,t}^{(1)\top} \mathbb{Y}_{gh,(t-1)}$, and we have

$$\begin{aligned} s_{31,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{\mathbb{Y}_{\cdot,h,t}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g^{(1)}, \cdot)})^\top \mathbb{Y}_{gh,(t-1)}\}, \\ \mathbf{s}_{32,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{(\mathbf{1}_{N_{2h}} \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_g^{(1)}, \cdot)})^\top \mathbb{Y}_{gh,(t-1)}\}. \end{aligned}$$

Since $\mathbb{Y}_{(t-1)}$ is independent with $\mathbf{X}_t^{(1)}$, we have $\mathbf{s}_{32,gh} = \mathbf{0}$. Note that

$$\begin{aligned} s_{31,gh} &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t E\{\mathbb{Y}_{\cdot,h,(t-1)}^\top (\mathbf{I}_{N_{2h}} \otimes \mathbf{W}_1^{(\mathcal{R}_g^{(1)}, \cdot)})^\top \mathbb{Y}_{gh,(t-1)}\} \\ &= \lim_{n,T \rightarrow \infty} (n^2 T)^{-1} \sum_t [\text{tr}(\mathcal{W}_{gh}^{(1)} \boldsymbol{\Sigma}_{\cdot,h,gh}^{\mathbf{X}}) + \text{tr}(\mathcal{W}_{gh}^{(1)} \boldsymbol{\Sigma}_{\cdot,h,gh}^{\mathcal{E}})] = \nu_{gh,gh}^1 + \kappa_{gh,gh}^1. \end{aligned}$$

Therefore, we obtain the final expression.

(4). Expression of \mathbf{s}_4 .

We next show that $\mathbf{s}_4 = (s_{4,gh}I(g = g', h = h') : \mathcal{I}_{gh} \in [G_1G_2], \mathcal{I}_{g'h'} \in [G_1G_2])$ with $s_{4,gh} = \nu_{gh,gh}^\alpha + \kappa_{gh,gh}^\alpha$. Since $\mathbf{M}^\alpha = (\mathbf{M}_{\mathcal{I}_{gh}\mathcal{I}_{g'h'}}^\alpha : \mathcal{I}_{gh} \in [G_1G_2], \mathcal{I}_{g'h'} \in [G_1G_2])$, where $\mathbf{M}_{\mathcal{I}_{gh}\mathcal{I}_{g'h'}}^\alpha = \sum_t \mathbb{Y}_{gh,(t-1)}^\top \mathbb{Y}_{gh,(t-1)} I(g = g', h = h')$, we derive $s_{4,gh} = \lim_{n,T \rightarrow \infty} (n^2T)^{-1} \sum_t E(\sum_t \mathbb{Y}_{gh,(t-1)}^\top \mathbb{Y}_{gh,(t-1)})$.

$$\begin{aligned} s_{4,gh} &= \lim_{n,T \rightarrow \infty} (n^2T)^{-1} \sum_t E(\mathbb{Y}_{gh,(t-1)}^\top \mathbb{Y}_{gh,(t-1)}) \\ &= \lim_{n,T \rightarrow \infty} (n^2T)^{-1} \sum_t [\text{tr}(\boldsymbol{\Sigma}_{gh,gh}^{\mathbf{X}}) + \text{tr}(\boldsymbol{\Sigma}_{gh,gh}^{\mathcal{E}})] = \nu_{gh,gh}^\alpha + \kappa_{gh,gh}^\alpha. \end{aligned}$$

Therefore, we obtain the final expression. ■

B.4 Additional Discussion of Weighted Least Squares Estimation with Group-specific Error Variances

B.4.1 TECHNICAL ASSUMPTION

We need the Assumption A.10 in Appendix B.4.1 instead of the Assumption 3 in Section 4.

Assumption A.10. (DISTRIBUTION OF NOISE TERM WITH GROUP-SPECIFIC VARIANCE) Assume $\varepsilon_{ij,t}$ is a zero-mean sub-Gaussian variable, and independently distributed across $i \in [N_1], j \in [N_2]$ and $t \in [T]$ and independent with $\{\mathbf{Y}_s : s \leq t-1\}$, $\{\mathbf{X}_s^{(1)} = (\mathbf{x}_{is}^{(1)} : i \in [N_1])^\top : s \leq t\}$ and $\{\mathbf{X}_s^{(2)} = (\mathbf{x}_{js}^{(2)} : j \in [N_2])^\top : s \leq t\}$. Assume that $\text{var}(\varepsilon_{ij,t}) = \sigma_{g_i^{(1)0} g_j^{(2)0}}^2$, where $g_i^{(1)0}$ and $g_j^{(2)0}$ are true group memberships. It holds $c_1 \leq \min_{g^{(1)}, g^{(2)}} \sigma_{g^{(1)}g^{(2)}}^2 \leq \max_{g^{(1)}, g^{(2)}} \sigma_{g^{(1)}g^{(2)}}^2 \leq c_2$, where c_1 and c_2 are two finite positive constants.

B.4.2 PROOF OF THEOREM 9 (I)

Proof We prove the proposition in three steps. In step (1), we first show that $d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) = O_p\{T^{-1}(\log(N_1N_2))^2\}$ holds under the clustered error case. This is a required condition for the Theorem 5. In step (2), we show that when $G_1 \geq G_{1,0}$, $G_2 \geq G_{2,0}$ and $c_{\text{gap}} \gg d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0)$, the three conclusions in Theorem 5 still hold. In the final step (3), we come to the proposition.

Step (1). Order of $d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0)$.

We first show that Lemma 17 and Lemma 23 holds in step (A), and then show that Theorem 2 holds in step (B).

STEP (A). LEMMA 17 AND LEMMA 23.

First, we show that (A.66) and (A.67) hold in Lemma 17. To validate (A.66), it suffices to show (A.68)–(A.70) hold. By scrutinizing the three inequality, we only need to take care of (A.68) and (A.69), which are related to the error term. Rewrite the two inequality by

$$P\left(\left|\frac{1}{T} \sum_t \varepsilon_{ij,t}^2 - \sigma_{g_i^{(1)} g_j^{(2)}}^2\right| > x/3\right) \leq 2 \exp\{-c_1 T \min(x^2, x)\}, \quad (\text{A.18})$$

$$P\left(\sup_{\|\boldsymbol{\xi}\|_{\max} < 2R} \left| \frac{2}{T} \sum_t \varepsilon_{ij,t} \mathcal{X}_{ij,t}^\top \boldsymbol{\xi} \right| > x/3\right) \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_3m\right\}, \quad (\text{A.19})$$

where (A.18) could be obtained by Lemma 27. Since $E(\varepsilon_{ij,t}) = 0$, (A.19) could be proved following the similar procedure to prove (A.69). Hence, we have

$$P\left(\sup_{\|\boldsymbol{\Theta}_{ij}\|_{\max} < R} \left| \frac{1}{T} Q_{ij}(\boldsymbol{\Theta}_{ij}) - \frac{1}{T} Q_{ij}^*(\boldsymbol{\Theta}_{ij}) \right| > x\right) \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_2m\right\}$$

still holds for clustered error variance. Likewise, (A.67) could be validated similarly.

Next, we show that Lemma 23 holds. First note that Lemma 33 still holds under Assumption 4 and Assumption 5, and the two assumptions could still be satisfied if the error terms have group-specific variance. Hence, (A.85)–(A.87) in Lemma 23 are obtained when Assumption 2 holds.

STEP (B). THEOREM 2 HOLDS.

By Step (A) we know that Lemma 17 and Lemma 23 hold, then we could borrow the idea of proof for Theorem 2, and show that $d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) = O_p\{T^{-1}(\log(N_1N_2))^2\}$.

Step (2). Proof of Theorem 5.

We prove the two statements in Lemma 25 first in Step (A)–(B), and then we prove the statement in Theorem 5 in Step (C).

STEP (A). PROOF OF STATEMENT (I) OF LEMMA 25.

Since the definition of $\mathcal{N}_\eta^{(1)}$ and $\mathcal{A}_\eta(\boldsymbol{\xi}, g^{(2)0}, \mathcal{G}_{-1}^0)$ are the same as the case under homogeneity, hence by the similar techniques in the proof of (i) of Theorem 5, one could keep the conclusion under heterogeneity case.

STEP (B). PROOF OF STATEMENT (II) OF LEMMA 25.

First, we show some required lemmas still hold. Note that in the clustered error case, the Assumption 2 for convexity and Assumption 6 still maintain. Since we have shown that Lemma 17 and Lemma 23 holds, we next show that Lemma 33 also holds. Because in the clustered error case, Assumption 4 and Assumption 5 are not affected, Lemma 33 could be proved similarly.

Then it comes to the proof for statement (ii). Similar as the proof procedure in Theorem 5, by using Lemma 23 and Assumption 2, one could have for any $g^{(2)} \notin \mathcal{A}_\eta(\boldsymbol{\theta}, g_j^{(2)0})$,

$$\frac{1}{N_1T} Q_j^*(\boldsymbol{\xi}_{g^{(2)}}^{(2)}; \boldsymbol{\xi}_{g^{(1)}}^{(1)}, \widehat{\mathcal{G}}_1(\boldsymbol{\xi})) - \frac{1}{N_1T} Q_j^*(\boldsymbol{\xi}_{g_j^{(2)0}}^{(2)}; \boldsymbol{\xi}_{g^{(1)0}}^{(1)0}, \mathcal{G}_1^0) \geq \tau_{\min}(c_{\text{gap}}c_\pi/2 - \eta), \quad (\text{A.20})$$

when $N_1 \rightarrow \infty$. On the other hand, for $\widetilde{g}_j^{(2)} \in \mathcal{A}_\eta^{(2)}(\boldsymbol{\xi}, g_j^{(2)0}, \mathcal{G}_1^0)$, the following still holds,

$$\frac{1}{N_1T} Q_j^*(\boldsymbol{\xi}_{\widetilde{g}_j^{(2)}}^{(2)}; \boldsymbol{\xi}_{g^{(1)}}^{(1)}, \widehat{\mathcal{G}}_1(\boldsymbol{\xi})) - \frac{1}{N_1T} Q_j^*(\boldsymbol{\xi}_{g_j^{(2)0}}^{(2)}; \boldsymbol{\xi}_{g^{(1)0}}^{(1)0}, \mathcal{G}_1^0) \leq \tau_{\max}\{d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) + \eta\}. \quad (\text{A.21})$$

Combining (A.20) and (A.21), and by the order of $d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0)$ derived in Step (1), similar conclusion as in the proof of Theorem 5 could be obtained,

$$W_{jg^{(2)}}(\boldsymbol{\xi}) \leq 2I \left(\sup_{i,j} \sup_{\|\boldsymbol{\Theta}_{ij}\|_{\max} < R} \left| \frac{1}{T} Q_{ij}(\boldsymbol{\Theta}_{ij}) - \frac{1}{T} Q_{ij}^*(\boldsymbol{\Theta}_{ij}) \right| \geq \epsilon_\eta/2 \right).$$

Further by Step (1), the results could be obtained.

STEP (C). PROOF OF STATEMENT (III) OF LEMMA 25.

We first show that Lemma 24 and Proposition 4 hold, and then (iii) could be proved. Since Theorem 2 holds, one could use the same procedure to prove Lemma 24. To prove Proposition 4, we need a new assumption for the distribution of noise term, namely Assumption A.10. Together with the previously stated Lemma 17, Lemma 23, Lemma 24 and Theorem 2, Proposition 4 holds under the Assumption 1–2, 4–5 and A.10.

Then it comes to the proof of (iii) of Lemma 25. Similar as the proof of (iii) for Theorem 5, first it could be obtained that

$$\begin{aligned} & \max_{\tilde{g}^{(2)} \in [G_2]} \min_{g^{(2)} \in [G_{2,0}]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(2)}}^c - \boldsymbol{\theta}_{g^{(2)}}^{(2)0}\|^2 + \frac{1}{N_1} \sum_i |\widehat{\alpha}_{\tilde{g}_i^{(1)}\tilde{g}^{(2)}} - \alpha_{g_i^{(1)0}g^{(2)}}^0|^2 \right) \\ & = O_p(c_\pi^{-1}T^{-1}(\log(N_1N_2))^2) \end{aligned}$$

by Proposition 4, and due to Lemma 24, it also holds that

$$\begin{aligned} & \max_{\tilde{g}^{(2)} \in [G_2]} \min_{\tilde{g}^{(2)} \in [G_2]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(2)}}^{(2)} - \boldsymbol{\theta}_{g^{(2)}}^{(2)0}\|^2 + \frac{1}{N_1} \sum_i |\widehat{\alpha}_{\tilde{g}_i^{(1)}\tilde{g}^{(2)}} - \alpha_{g_i^{(1)0}g^{(2)}}^0|^2 \right) \\ & = O_p(c_\pi^{-1}T^{-1}(\log(N_1N_2))^2). \end{aligned}$$

Then we have $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(2)}$ with $\eta = \tau_{\min}c_{\text{gap}}c_\pi/\{8(\tau_{\min} + \tau_{\max})\}$ under the condition that $c_{\text{gap}} \gg T^{-1}\log(N_1N_2)^2$. Then the statements (i) and (ii) in Lemma 25 hold, and then one could demonstrate Theorem 5 by using the conclusion Lemma 25. Other proof details are consistent with those in the i.i.d. error case.

Step (3). Membership Consistency.

Similarly as the proof for Corollary 6, using the statement in Theorem 5 proved in Step (2), the arguments (34) and (35) hold. \blacksquare

B.4.3 PROOF OF THEOREM 9 (II)

Proof We prove the result by two steps. In step (1), we first show that $|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2} - \sigma_{g^{(1)}g^{(2)}}^{-2}| = o_p(1)$ when $\log(N_{1g^{(1)}}N_{2g^{(2)}})^2/T \rightarrow 0$. Then in step (2), we show that $\sup_{g^{(1)}, g^{(2)}} |\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2} - \sigma_{g^{(1)}g^{(2)}}^{-2}| = o_p(1)$ under the condition $\sqrt{T} \gg \log(G_1G_2n^2) + m$.

Step (1). Order of $|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2} - \sigma_{g^{(1)}g^{(2)}}^{-2}|$.

By Assumption A.10, and define the event $\mathcal{O}_{g^{(1)}g^{(2)}} = \{\widehat{\sigma}_{g^{(1)}g^{(2)}}^2 \geq \sigma_{g^{(1)}g^{(2)}}^2/2\} = \{\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2} \leq 2\sigma_{g^{(1)}g^{(2)}}^{-2} \leq 2c\}$, where c is a positive constant. Then, we have

$$\begin{aligned} P\left\{|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2} - \sigma_{g^{(1)}g^{(2)}}^{-2}| \geq t\right\} &= P\left\{\left|\sigma_{g^{(1)}g^{(2)}}^{-2}(\sigma_{g^{(1)}g^{(2)}}^2 - \widehat{\sigma}_{g^{(1)}g^{(2)}}^2)\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2}\right| \geq t\right\} \\ &\leq P\left\{\left|\sigma_{g^{(1)}g^{(2)}}^{-4}(\sigma_{g^{(1)}g^{(2)}}^2 - \widehat{\sigma}_{g^{(1)}g^{(2)}}^2)\right| \geq t/2\right\} + P(\mathcal{O}_{g^{(1)}g^{(2)}}^c) \\ &\leq P\left\{\left|\sigma_{g^{(1)}g^{(2)}}^2 - \widehat{\sigma}_{g^{(1)}g^{(2)}}^2\right| \geq \frac{t}{2c^2}\right\} + P(\mathcal{O}_{g^{(1)}g^{(2)}}^c). \end{aligned}$$

Note that for the event $\mathcal{O}_{g^{(1)g^{(2)}}}^c$, we could calculate the probability as

$$\begin{aligned} P(\mathcal{O}_{g^{(1)g^{(2)}}}^c) &= P\left\{\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 < \frac{1}{2}\sigma_{g^{(1)g^{(2)}}}^2\right\} = P\left\{\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2 \leq -\frac{1}{2}\sigma_{g^{(1)g^{(2)}}}^2\right\} \\ &\leq P\left\{|\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2| \geq \frac{1}{2}\sigma_{g^{(1)g^{(2)}}}^2\right\} \\ &\leq P\left\{|\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2| \geq \frac{1}{2c}\right\}. \end{aligned}$$

Hence, we have

$$\begin{aligned} &P\left\{|\widehat{\sigma}_{g^{(1)g^{(2)}}}^{-2} - \sigma_{g^{(1)g^{(2)}}}^{-2}| \geq t\right\} \\ &\leq P\left\{\left|\sigma_{g^{(1)g^{(2)}}}^2 - \widehat{\sigma}_{g^{(1)g^{(2)}}}^2\right| \geq \frac{t}{2c^2}\right\} + P\left\{|\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2| \geq \frac{1}{2c}\right\} \\ &\leq P\left\{\left|\sigma_{g^{(1)g^{(2)}}}^2 - \widehat{\sigma}_{g^{(1)g^{(2)}}}^2\right| \geq \frac{t}{2c^2}\right\} + P\left\{|\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2| \geq \frac{t}{2c}\right\} \\ &\leq 2P\left\{\left|\sigma_{g^{(1)g^{(2)}}}^2 - \widehat{\sigma}_{g^{(1)g^{(2)}}}^2\right| \geq \frac{t}{C}\right\}. \end{aligned}$$

when $0 < t < 1$ and $C = \max(2c^2, 2c)$. Note that

$$\begin{aligned} |\widehat{\sigma}_{g^{(1)g^{(2)}}}^2 - \sigma_{g^{(1)g^{(2)}}}^2| &= (N_{1g^{(1)}}N_{2g^{(2)}}T)^{-1} \left| \sum_{i \in \mathcal{R}_{g^{(1)}}^{(1)}} \sum_{j \in \mathcal{R}_{g^{(2)}}^{(2)}} \{Q_{ij}(\widehat{\Theta}_{ij}) - Q_{ij}^*(\Theta_{ij}^0)\} \right| \\ &\leq (N_{1g^{(1)}}N_{2g^{(2)}}T)^{-1} \left[\left| \sum_{i \in \mathcal{R}_{g^{(1)}}^{(1)}} \sum_{j \in \mathcal{R}_{g^{(2)}}^{(2)}} \{Q_{ij}(\widehat{\Theta}_{ij}) - Q_{ij}^*(\widehat{\Theta}_{ij})\} \right| \right. \\ &\quad \left. + \left| \sum_{i \in \mathcal{R}_{g^{(1)}}^{(1)}} \sum_{j \in \mathcal{R}_{g^{(2)}}^{(2)}} \{Q_{ij}^*(\widehat{\Theta}_{ij}) - Q_{ij}^*(\Theta_{ij}^0)\} \right| \right] \end{aligned}$$

Then we derive the order of the two terms on the right hands separately. First, by Lemma 17, we could directly obtain that

$$\begin{aligned} &(N_{1g^{(1)}}N_{2g^{(2)}}T)^{-1} \left| \sum_{i \in \mathcal{R}_{g^{(1)}}^{(1)}} \sum_{j \in \mathcal{R}_{g^{(2)}}^{(2)}} \{Q_{ij}(\widehat{\Theta}_{ij}) - Q_{ij}^*(\widehat{\Theta}_{ij})\} \right| \\ &\leq (N_{1g^{(1)}}N_{2g^{(2)}}T)^{-1} \left| \sum_{i \in \mathcal{R}_{g^{(1)}}^{(1)}} \sum_{j \in \mathcal{R}_{g^{(2)}}^{(2)}} \sup_{\|\Theta\|_{\max} \leq R} \{Q_{ij}(\Theta) - Q_{ij}^*(\Theta)\} \right| \\ &\leq \sup_{i \in \mathcal{R}_{g^{(1)}}^{(1)}, j \in \mathcal{R}_{g^{(2)}}^{(2)}} \sup_{\|\Theta\|_{\max} \leq R} |T^{-1} \{Q_{ij}(\Theta) - Q_{ij}^*(\Theta)\}|, \end{aligned}$$

and

$$\begin{aligned} &P\left\{\sup_{i \in \mathcal{R}_{g^{(1)}}^{(1)}, j \in \mathcal{R}_{g^{(2)}}^{(2)}} \sup_{\|\Theta\|_{\max} \leq R} |T^{-1} \{Q_{ij}(\Theta) - Q_{ij}^*(\Theta)\}| \geq t\right\} \\ &\leq N_{1g^{(1)}}N_{2g^{(2)}} \exp\{-c_1 \min(Tt^2, \sqrt{T}t) + c_2 m\}. \end{aligned}$$

Then, for the second term, one could borrow the idea of the proof of Lemma 23 to obtain,

$$\begin{aligned} & P\left\{\frac{1}{N_{1g^{(1)}}N_{2g^{(2)}}T}\left|\sum_{i\in\mathcal{R}_{g^{(1)}}^{(1)}}\sum_{j\in\mathcal{R}_{g^{(2)}}^{(2)}}\{Q_{ij}^*(\widehat{\Theta}_{ij})-Q_{ij}^*(\Theta_{ij}^0)\}\right|\geq t\right\} \\ & \leq N_{1g^{(1)}}N_{2g^{(2)}}\exp\{-c_1\min(Tt^2,\sqrt{T}t)+c_2m\}. \end{aligned}$$

Hence, we have

$$\begin{aligned} P\left\{|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2}-\sigma_{g^{(1)}g^{(2)}}^{-2}|\geq t\right\} & \leq 2P\left\{\left|\sigma_{g^{(1)}g^{(2)}}^2-\widehat{\sigma}_{g^{(1)}g^{(2)}}^2\right|\geq\frac{t}{C}\right\} \\ & \leq N_{1g^{(1)}}N_{2g^{(2)}}\exp\{-c_1\min(Tt^2,\sqrt{T}t)+c_2m\}. \end{aligned}$$

Step (2). Order of $\sup_{g^{(1)}g^{(2)}}|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2}-\sigma_{g^{(1)}g^{(2)}}^{-2}|$.

By step (1), we have

$$P\left\{\sup_{g^{(1)}g^{(2)}}|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2}-\sigma_{g^{(1)}g^{(2)}}^{-2}|\geq t\right\}\leq G_1G_2n^2\exp\{-c_1\min(Tt^2,\sqrt{T}t)+c_2m\},$$

where n is defined in the statement of Theorem 9. This implies that $\sup_{g^{(1)}g^{(2)}}|\widehat{\sigma}_{g^{(1)}g^{(2)}}^{-2}-\sigma_{g^{(1)}g^{(2)}}^{-2}|=O_p\{(\log(G_1G_2n^2)+m)/\sqrt{T}\}$, and thus the proof is completed. \blacksquare

B.4.4 PROOF OF THEOREM 9 (III)

Proof Note that

$$nT^{1/2}(\widehat{\boldsymbol{\theta}}^w-\boldsymbol{\theta}^0)=nT^{1/2}\widetilde{\mathbf{M}}^{-1}\check{\boldsymbol{\delta}}=((n^2T)^{-1}\widetilde{\mathbf{M}})^{-1}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}}\stackrel{\text{def}}{=} \widetilde{\mathbf{M}}_{nT}^{-1}\boldsymbol{\Lambda}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}},$$

where $\check{\boldsymbol{\delta}}=(\check{\boldsymbol{\delta}}^{(1)\top},\check{\boldsymbol{\delta}}^{(2)\top},\check{\boldsymbol{\delta}}^{\alpha\top})^\top$ and

$$\check{\boldsymbol{\delta}}_{g^{(1)}}^{(1)}=\sum_{t,g^{(2)}}(\widehat{\sigma}_{g^{(1)}g^{(2)}})^{-2}\mathbb{X}_{g^{(1)}g^{(2)},t}^{(1)\top}\mathbb{E}_{g^{(1)}g^{(2)},t},\quad \check{\boldsymbol{\delta}}^{(1)}=(\check{\boldsymbol{\delta}}_{g^{(1)}}^{(1)\top}:g^{(1)}\in[G_1])^\top,$$

$$\check{\boldsymbol{\delta}}_{g^{(2)}}^{(2)}=\sum_{t,g^{(1)}}(\widehat{\sigma}_{g^{(1)}g^{(2)}})^{-2}\mathbb{X}_{g^{(1)}g^{(2)},t}^{(2)\top}\mathbb{E}_{g^{(1)}g^{(2)},t},\quad \check{\boldsymbol{\delta}}^{(2)}=(\check{\boldsymbol{\delta}}_{g^{(2)}}^{(2)\top}:g^{(2)}\in[G_2])^\top,$$

$$\check{\boldsymbol{\delta}}_{g^{(1)}g^{(2)}}^\alpha=\sum_t(\widehat{\sigma}_{g^{(1)}g^{(2)}})^{-2}\mathbb{Y}_{g^{(1)}g^{(2)}(t-1)}^\top\mathbb{E}_{g^{(1)}g^{(2)},t},\quad \check{\boldsymbol{\delta}}^\alpha=(\check{\boldsymbol{\delta}}_{g^{(1)}g^{(2)}}^{\alpha\top}:g^{(1)}g^{(2)}\in[g^{(1)}g^{(2)}])^\top.$$

Then, we have

$$\begin{aligned} nT^{1/2}\boldsymbol{\eta}^\top(\widehat{\boldsymbol{\theta}}^w-\boldsymbol{\theta}^0) & = \boldsymbol{\eta}^\top\widetilde{\mathbf{M}}_{nT}^{-1}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}} \\ & = \boldsymbol{\eta}^\top\{\widetilde{\mathbf{M}}_{nT}^{-1}-(\widetilde{\mathbf{M}}_{nT}^0)^{-1}\}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}}+\boldsymbol{\eta}^\top(\widetilde{\mathbf{M}}_{nT}^0)^{-1}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}} \\ & = \boldsymbol{\eta}^\top(\widetilde{\mathbf{M}}_{nT}^0)^{-1}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}}+\boldsymbol{\eta}^\top\widetilde{\mathbf{M}}_{nT}^{-1}\{\widetilde{\mathbf{M}}_{nT}^0-\widetilde{\mathbf{M}}_{nT}\}(\widetilde{\mathbf{M}}_{nT}^0)^{-1}(nT^{1/2})^{-1}\check{\boldsymbol{\delta}} \end{aligned}$$

$$\begin{aligned}
 &= \boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} \widetilde{\boldsymbol{\delta}} + \boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}}) \\
 &+ \boldsymbol{\eta}^\top \widetilde{\mathbf{M}}_{nT}^{-1} \{ \widetilde{\mathbf{M}}_{nT}^0 - \widetilde{\mathbf{M}}_{nT} \} (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} \check{\boldsymbol{\delta}}.
 \end{aligned}$$

Hence, it suffices to show that

$$\boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} \widetilde{\boldsymbol{\delta}} \rightarrow_d N(0, \boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}^0)^{-1} \boldsymbol{\eta}), \quad (\text{A.22})$$

$$\boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}}) = o_p(1), \quad (\text{A.23})$$

$$\boldsymbol{\eta}^\top \widetilde{\mathbf{M}}_{nT}^{-1} \{ \widetilde{\mathbf{M}}_{nT}^0 - \widetilde{\mathbf{M}}_{nT} \} (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} \check{\boldsymbol{\delta}} = o_p(1). \quad (\text{A.24})$$

While it is obvious that (A.24) holds when $n, T \rightarrow \infty$, we prove the (A.22)–(A.23) in the following two steps.

1. Proof of (A.22)

By noting that (A.22) and (A.56) have similar forms, one could borrow the idea from the proof for (A.56) with the new Assumption A.10. Then we could obtain the similar results by using the central limit theorem of the martingale difference array.

2. Proof of (A.23)

Note that $|\boldsymbol{\eta}^\top (\widetilde{\mathbf{M}}_{nT}^0)^{-1} (nT^{1/2})^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}})| \leq \|(\widetilde{\mathbf{M}}_{nT}^0)^{-1} \boldsymbol{\eta}\| \|(nT^{1/2})^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}})\|$, and we have assumed that $\lambda_{\min}(\mathbf{M}^0) > 0$ and $\max_{g^{(1)}, g^{(2)}} \sigma_{g^{(1)}, g^{(2)}}^2 < \infty$, hence it holds that $\|(\widetilde{\mathbf{M}}_{nT}^0)^{-1} \boldsymbol{\eta}\| < c < \infty$. Then to derive the order of $\|(nT^{1/2})^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}})\|$, we focus on the order of $(n^2T)^{-1} (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}})^\top (\check{\boldsymbol{\delta}} - \widetilde{\boldsymbol{\delta}})$. We take the term $\check{\boldsymbol{\delta}}^{(1)} \in \mathbb{R}^{G \times (p_1+1)}$ and $\widetilde{\boldsymbol{\delta}}^{(1)} \in \mathbb{R}^{G \times (p_1+1)}$ for example to illustrate. Note that

$$\begin{aligned}
 &(n^2T)^{-1} (\check{\boldsymbol{\delta}}^{(1)} - \widetilde{\boldsymbol{\delta}}^{(1)})^\top (\check{\boldsymbol{\delta}}^{(1)} - \widetilde{\boldsymbol{\delta}}^{(1)}) \\
 &= (n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} (\widehat{\sigma}_{g^{(1)}, g^{(2)}}^{-2} - \sigma_{g^{(1)}, g^{(2)}}^{-2})^2 \mathbb{E}_{g^{(1)}, g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top} \mathbb{E}_{g^{(1)}, g^{(2)}, t} \\
 &\leq (n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} \left\{ \sup_{g^{(1)}, g^{(2)}} (\widehat{\sigma}_{g^{(1)}, g^{(2)}}^{-2} - \sigma_{g^{(1)}, g^{(2)}}^{-2})^2 \right\} \mathbb{E}_{g^{(1)}, g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top} \mathbb{E}_{g^{(1)}, g^{(2)}, t}.
 \end{aligned}$$

By Theorem 9 (ii), one has $\sup_{g^{(1)}, g^{(2)}} (\widehat{\sigma}_{g^{(1)}, g^{(2)}}^{-2} - \sigma_{g^{(1)}, g^{(2)}}^{-2})^2 = O_p\{T^{-1}(m + \log(G_1 G_2 n^2))\}$, hence we next derive the order of $(n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} E(\mathbb{E}_{g^{(1)}, g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top} \mathbb{E}_{g^{(1)}, g^{(2)}, t})$.

By Assumption A.10 and $\max_{g^{(1)}, g^{(2)}} \sigma_{g^{(1)}, g^{(2)}} \leq c$, we know that with $\boldsymbol{\eta}_1 \in \mathbb{R}^{p+1}$ and $\|\boldsymbol{\eta}_1\| = 1$,

$$\begin{aligned}
 &(n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} E(\mathbb{E}_{g^{(1)}, g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top} \mathbb{E}_{g^{(1)}, g^{(2)}, t}) \\
 &= (n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} E\{\sigma_{g^{(1)}, g^{(2)}}^2 \text{tr}(\mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top})\} \\
 &\leq c^2 (n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} E\{\text{tr}(\mathbf{W}_1^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot)} \mathbf{Y}_t^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})} \mathbf{Y}_t^{(\cdot, \mathcal{R}_{g^{(2)}}^{(2)})\top} \mathbf{W}_1^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)}, \cdot)\top})\}
 \end{aligned}$$

$$\begin{aligned}
 & + c^2(n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} E\{\text{vec}(\mathbf{1}_{N_{2g^{(2)}}}) \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})^\top} \text{vec}(\mathbf{1}_{N_{2g^{(2)}}}) \otimes (\mathbf{X}_t^{(1)})^{(\mathcal{R}_{g^{(1)}, \cdot}^{(1)})}\} \\
 & = c^2(n^2T)^{-1} \sum_t E\{\text{tr}(\mathbf{W}_1 \mathbf{Y}_t \mathbf{Y}_t^\top \mathbf{W}_1^\top)\} + c^2 N_2 (n^2T)^{-1} \sum_t \sum_{i \in [N_1]} E\|\mathbf{x}_{it}^{(1)}\|^2.
 \end{aligned}$$

Note that

$$\sum_t E\{\text{tr}(\mathbf{W}_1 \mathbf{Y}_t \mathbf{Y}_t^\top \mathbf{W}_1^\top)\} \leq c_1 \sum_t \text{tr}(\mathbf{W}_1 \mathbf{1}_{N_1} \mathbf{1}_{N_2}^\top \mathbf{1}_{N_2} \mathbf{1}_{N_1}^\top \mathbf{W}_1^\top) = O(n^2T),$$

where $c_1 = E(Y_{ij,t}^2) < \infty$ by Lemma 32, and also note that

$$\sum_t \sum_{i \in [N_1]} E\|\mathbf{x}_{it}^{(1)}\|^2 \leq c_{Kp_1} T N_1 = O(nTs)$$

due to the K -convexity property of $\mathbf{x}_{it}^{(1)}$ and the definition of n , we have that

$$(n^2T)^{-1} \sum_{g^{(1)}, g^{(2)}, t} (\mathbb{E}_{g^{(1)}, g^{(2)}, t}^\top \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)} \mathbb{X}_{g^{(1)}, g^{(2)}, t}^{(1)\top} \mathbb{E}_{g^{(1)}, g^{(2)}, t}) = O_p(1).$$

Together with $\sup_{g^{(1)}, g^{(2)}} (\widehat{\sigma}_{g^{(1)}, g^{(2)}}^{-2} - \sigma_{g^{(1)}, g^{(2)}}^{-2})^2 = O_p\{T^{-1}(m + \log(GHn^2))\}$, we know that

$$(n^2T)^{-1} (\check{\boldsymbol{\delta}}^{(1)} - \tilde{\boldsymbol{\delta}}^{(1)})^\top (\check{\boldsymbol{\delta}}^{(1)} - \tilde{\boldsymbol{\delta}}^{(1)}) = O_p\{(m + \log(G_1 G_2 n^2))/T\} = o_p(1)$$

when $T \gg (\log(G_1 G_2 n^2))$ due to that m is a fixed constant.

By step 1–2, we could arrive at the final conclusion. ■

B.4.5 SIMULATION RESULTS OF WEIGHTED LEAST SQUARES ESTIMATION

In this section, we conduct the simulation experiment when $G_{1,0} = G_{2,0} = 3$. We set the group-specific error variances as $\tilde{\boldsymbol{\Sigma}} = (\sigma_{g^{(1)}, g^{(2)}} : g^{(1)} \in [G_{1,0}], g^{(2)} \in [G_{2,0}]) \in \mathbb{R}^{3 \times 3}$ as

$$\tilde{\boldsymbol{\Sigma}} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}.$$

Then we set $(N_1, N_2) \in \{(100, 80), (200, 150)\}$ and set $T \in \{20, 40, 80\}$. Next, we repeat the experiment for $R = 300$ times, and show the RMSEs results of $\tilde{\boldsymbol{\theta}}^w$ in Table A.5. From the table, we first focus on the estimation results under weighted least squares method (marked as ‘‘WLS’’ in the table). We could see that when the error has group-specific variances, our proposed estimation could estimated the true group memberships with high accuracy as the sample size grows. Furthermore, we note that RMSEs of all parameters decrease toward the corresponding oracle results when the sample sizes are large. For the inference results of WLS, when both N_1, N_2, T are large, the CPs are all around 0.95. Then, we compare the

results of WLS and the ordinary method in the main text (marked as ‘‘OLS’’). One could see that the RMSEs of WLS estimator are always smaller than those of OLS estimator, and the RMSEs of the corresponding oracle estimators have the same pattern. This shows that when the error has group-specific variances, our proposed weighted least squares estimator outperforms the ordinary least squares estimator.

Appendix C. Specific Expressions of \mathbf{M} in Equation (13)

$$\mathbf{M}_{g^{(l)}}^{(l)} = \sum_{t, g^{-(l)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)}, \quad \mathbf{M}^{(l)} = \text{diag}(\mathbf{M}_{g^{(l)}}^{(l)} : g^{(l)} \in [G_l]) \in \mathbb{R}^{G_l(p_l+1) \times G_l(p_l+1)},$$

$$\mathbf{M}_{g^{(l)}g^{(m)}}^{(lm)} = \sum_{t, g^{-(l,m)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(m)},$$

$$\mathbf{M}^{(lm)} = (\mathbf{M}_{g^{(l)}g^{(m)}}^{(lm)} : g^{(l)} \in [G_l], g^{(m)} \in [G_m]) \in \mathbb{R}^{G_l(p_l+1) \times G_m(p_m+1)},$$

$$\mathbf{M}_{g^{(l)}\mathcal{I}_{g^{(l)'}, g^{-(l)}}}^{(l\alpha)} = \sum_t \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)} I(g^{(l)} = g^{(l)'}),$$

$$\text{with } \mathcal{I}_{g^{(l)'}, g^{-(l)}} = g^{(l)'} + \sum_{m \neq l} ((g^{(m)} - 1)G_l),$$

$$\mathbf{M}^{(l\alpha)} = \left(\mathbf{M}_{g^{(l)}\mathcal{I}_{g^{(l)'}, g^{-(l)}}}^{(l\alpha)} : g^{(l)} \in [G_l], \mathcal{I}_{g^{(l)'}, g^{-(l)}} \in \left[\prod_l G_l \right] \right) \in \mathbb{R}^{G_l(p_l+1) \times (\prod_l G_l)},$$

$$\mathbf{M}_{\mathcal{I}_{g^{(1)}, \dots, g^{(q)}}, \mathcal{I}_{g^{(1)'}, \dots, g^{(q)'}}}^{\alpha} = \sum_t \|\mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}\|^2 I\{(g^{(1)}, \dots, g^{(q)}) = (g^{(1)'}, \dots, g^{(q)'})\},$$

$$\mathbf{M}^{\alpha} = \left(\mathbf{M}_{\mathcal{I}_{g^{(1)}, \dots, g^{(q)}}, \mathcal{I}_{g^{(1)'}, \dots, g^{(q)'}}}^{\alpha} : \mathcal{I}_{g^{(1)}, \dots, g^{(q)}} \in \left[\prod_l G_l \right], \mathcal{I}_{g^{(1)'}, \dots, g^{(q)'}} \in \left[\prod_l G_l \right] \right)$$

$$\in \mathbb{R}^{(\prod_l G_l) \times (\prod_l G_l)},$$

$$\mathbf{b}_{g^{(l)}}^{(l)} = \sum_{t, g^{-(l)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, t}, \quad \mathbf{b}^{(l)} = (\mathbf{b}_{g^{(l)}}^{(l)} : g^{(l)} \in [G_l]) \in \mathbb{R}^{G_l(p_l+1)},$$

$$\mathbf{b}_{\mathcal{I}_{g^{(1)}, \dots, g^{(q)}}}^{\alpha} = \sum_t \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^{\top} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, t},$$

$$\mathbf{b}^{\alpha} = \left(\mathbf{b}_{\mathcal{I}_{g^{(1)}, \dots, g^{(q)}}}^{\alpha} : \mathcal{I}_{g^{(1)}, \dots, g^{(q)}} \in \left[\prod_l G_l \right] \right) \in \mathbb{R}^{\prod_l G_l},$$

where $\sum_{t, g^{-(l,m)}}$ is the simplified notation for $\sum_{\{g^{(1)}, \dots, g^{(q)}\} \setminus \{g^{(l)}, g^{(m)}\}}$.

Appendix D. GTNAR Algorithm

We summarize the GTNAR algorithm in the following Algorithm A.2.

Appendix E. Local Convergence of the Numerical Algorithm

In this section, we would show that our estimation Algorithm A.2 is convergent, and its solution is indeed a local minimizer of (9). We first show the convergence of the algorithm in (1), followed by the local minimization proof in (2).

Table A.5: RMSEs ($\times 1000$) of estimated parameters under scenario 3 ($G_{1,0} = 3, G_{2,0} = 3$) with 300 replications. The performances are evaluated for different sample sizes N_1, N_2 and the time length T . ‘‘OLS’’ means the results estimated using the ordinary least squares method proposed in Section 3, ‘‘WLS’’ means the results using weighted least squares proposed in Appendix B.4. The corresponding CPs are shown in the parenthesis.

N_1	N_2	T	Estimation	$\hat{\chi}^{(1)}$	$\hat{\chi}^{(2)}$	$\hat{\zeta}^{(1)}$	$\hat{\zeta}^{(2)}$	$\hat{\alpha}$	$\hat{\chi}^{(1)or}$	$\hat{\chi}^{(2)or}$	$\hat{\zeta}^{(1)or}$	$\hat{\zeta}^{(2)or}$	$\hat{\alpha}^{or}$	η_1	η_2		
100	80	40	OLS	24.79 (0.71)	12.47 (0.919)	87.59 (0.774)	52.22 (0.93)	46.10 (0.696)	15.10 (0.883)	11.18 (0.93)	51.11 (0.945)	50.26 (0.936)	14.68 (0.936)	0.1960	0.0077		
			WLS	24.45 (0.673)	12.02 (0.928)	68.28 (0.576)	49.96 (0.951)	45.45 (0.635)	11.20 (0.952)	28.77 (0.947)	47.34 (0.954)	14.53 (0.955)					
		OLS	15.86 (0.774)	8.31 (0.949)	53.54 (0.87)	33.93 (0.952)	22.80 (0.861)	9.39 (0.909)	8.30 (0.952)	39.10 (0.952)	33.94 (0.952)	10.71 (0.95)			0.0877	0.0000	
	20	80	WLS	14.50 (0.783)	7.61 (0.949)	39.13 (0.765)	32.16 (0.947)	22.56 (0.81)	7.24 (0.948)	7.38 (0.959)	22.81 (0.946)	31.87 (0.946)	10.63 (0.949)				
			OLS	23.50 (0.584)	12.61 (0.818)	87.96 (0.616)	41.33 (0.904)	52.51 (0.562)	11.6 (0.876)	8.55 (0.912)	36.19 (0.947)	38.21 (0.929)	10.66 (0.924)			0.2564	0.0189
		WLS	25.15 (0.524)	12.5 (0.827)	79.29 (0.409)	40.03 (0.47)	51.98 (0.47)	8.47 (0.95)	8.09 (0.933)	20.16 (0.957)	36.66 (0.95)	10.63 (0.954)					
200	150	40	OLS	12.26 (0.706)	5.48 (0.928)	40.14 (0.84)	26.10 (0.929)	16.32 (0.79)	7.85 (0.867)	5.48 (0.93)	26.85 (0.945)	26.10 (0.928)	7.59 (0.933)	0.1033	0.0000		
			WLS	11.06 (0.946)	5.09 (0.942)	32.77 (0.673)	25.12 (0.949)	16.07 (0.74)	5.66 (0.95)	4.98 (0.94)	15.58 (0.952)	24.89 (0.949)	7.54 (0.95)				
		OLS	5.89 (0.863)	4.21 (0.910)	17.95 (0.945)	19.14 (0.919)	5.49 (0.919)	5.80 (0.864)	4.21 (0.910)	17.71 (0.949)	19.14 (0.920)	5.29 (0.922)			0.0078	0.0000	
	80	40	WLS	4.12 (0.950)	3.90 (0.949)	10.83 (0.939)	18.19 (0.948)	5.42 (0.944)	3.99 (0.952)	3.89 (0.944)	10.46 (0.954)	18.17 (0.948)	5.23 (0.951)				
			OLS														
		WLS															

Algorithm A.2 Estimation of the GTNAR Model

- 1: **Input:** $\{\mathcal{Y}_t, \mathbf{X}_t^{(l)}, \mathbf{W}^{(l)}, G_l\}$ for $1 \leq l \leq q$.
- 2: Obtain initial group memberships $\mathcal{G}_l^{[0]}$ according to Algorithm A.3 in Appendix F. Let $\{\boldsymbol{\xi}^{[k]}, \mathcal{G}_l^{[k]}\}$ be the estimators and memberships in the k th iteration.
- 3: Repeat STEP 1 and STEP 2 for $k = 1, 2, \dots$ until convergence.
 - STEP 1. Given $\{\mathcal{G}_l^{[k-1]}\}$ for all layers $l \in [q]$, calculate $\boldsymbol{\xi}^{[k-1]} = (\mathbf{M}^{[k-1]})^{-1} \mathbf{b}^{[k-1]}$, where $\mathbf{M}^{[k-1]}$ and $\mathbf{b}^{[k-1]}$ are obtained from (13) with $\mathcal{G}_l^{[k-1]}$ s specified.
 - STEP 2. Given $\boldsymbol{\xi}^{[k-1]}$, sequentially update memberships $\mathcal{G}_l^{[k]} = (g_{i_l}^{(l)[k]} : 1 \leq i_l \leq N_l)^\top$ for $1 \leq l \leq q$ as follows,

$$\begin{aligned}
 g_{i_l}^{(l)[k]} = & \arg \min_{g_{i_l}^{(l)} \in [G_l]} \sum_{i_{-l}} \sum_{t=1}^T \left\{ Y_{i_1 i_2 \dots i_q, t} - \left(\sum_{l'=1}^{l-1} \lambda_{g_{i_{l'}}^{(l')[k]}} \sum_{m=1}^{N_{l'}} \frac{a_{i_{l'} m}^{(l')}}{n_{l' i_{l'}}} Y_{i_1 \dots i_{l'-1} m i_{l'+1} \dots i_q, (t-1)} \right. \right. \\
 & + \lambda_{g_{i_l}^{(l)[k]}} \sum_{m=1}^{N_l} \frac{a_{i_l m}^{(l)}}{n_{l i_l}} Y_{i_1 \dots i_{l-1} m i_{l+1} \dots i_q, (t-1)} + \sum_{l''=l+1}^q \lambda_{g_{i_{l''}}^{(l'')[k-1]}} \sum_{m=1}^{N_{l''}} \frac{a_{i_{l''} m}^{(l'')}}{n_{l'' i_{l''}}} Y_{i_1 \dots i_{l''-1} m i_{l''+1} \dots i_q, (t-1) \Big\} \\
 & - \alpha_{g_{i_1}^{(1)[k]} \dots g_{i_{l-1}}^{(l-1)[k]} g_{i_l}^{(l)} g_{i_{l+1}}^{(l+1)[k-1]} \dots g_{i_q}^{(q)[k-1]} Y_{i_1 i_2 \dots i_q, (t-1)} \\
 & - \left(\sum_{l'=1}^{l-1} \mathbf{x}_{i_{l'} t}^{(l')\top} \boldsymbol{\zeta}_{g_{i_{l'}}^{(l')[k]}}^{(l')[k]} + \mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}^{(l)[k]}}^{(l)[k]} + \sum_{l''=l+1}^q \mathbf{x}_{i_{l''} t}^{(l'')\top} \boldsymbol{\zeta}_{g_{i_{l''}}^{(l'')[k-1]}}^{(l'')[k-1]} \right) \Big\}^2 \tag{A.25}
 \end{aligned}$$

- 4: **Output:** Final estimator and memberships: $\widehat{\boldsymbol{\xi}} = \boldsymbol{\xi}^{[K]}$ and $\widehat{\mathcal{G}} = \{\mathcal{G}_l^{[K]} : l \in [q]\}$. Here K is the final number of iteration rounds.
-

(1) Convergence of the algorithm.

We show that the objective function defined in (9) is monotonically decreasing and hence could obtain the convergence. Denote the estimators in the k th iteration as $(\boldsymbol{\xi}^{(k)}, \mathcal{G}^{(k)})$. Then we have that

$$\begin{aligned}
 Q(\boldsymbol{\xi}^{(k+1)}, \mathcal{G}^{(k+1)}) &= \min_{\boldsymbol{\xi}} Q(\boldsymbol{\xi}, \mathcal{G}^{(k+1)}) \leq Q(\boldsymbol{\xi}^{(k)}, \mathcal{G}^{(k+1)}) \\
 &\stackrel{\textcircled{1}}{\leq} Q(\boldsymbol{\xi}^{(k)}, \{\mathcal{G}_1^{(k)}, \mathcal{G}_2^{(k+1)}, \dots, \mathcal{G}_q^{(k+1)}\}) \\
 &\stackrel{\textcircled{2}}{\leq} \dots\dots \\
 &\stackrel{\textcircled{3}}{\leq} Q(\boldsymbol{\xi}^{(k)}, \{\mathcal{G}_1^{(k)}, \mathcal{G}_2^{(k)}, \dots, \mathcal{G}_q^{(k)}\}) = Q(\boldsymbol{\xi}^{(k)}, \mathcal{G}^{(k)}),
 \end{aligned}$$

where the inequalities ①–③ are derived by

$$\mathcal{G}_l^{(k+1)} = \operatorname{argmin}_{\mathcal{G}_l} Q(\boldsymbol{\xi}^{(k)}, \{\mathcal{G}_1^{(k+1)}, \dots, \mathcal{G}_{q-1}^{(k+1)}, \mathcal{G}_l, \mathcal{G}_{l+1}^{(k)}, \dots, \mathcal{G}_q^{(k)}\})$$

in our updating mechanism (14). Then, we have that

$$Q(\boldsymbol{\xi}^{(0)}, \mathcal{G}^{(0)}) \geq Q(\boldsymbol{\xi}^{(1)}, \mathcal{G}^{(1)}) \geq \dots \geq Q(\boldsymbol{\xi}^{(k^*)}, \mathcal{G}^{(k^*)}) = Q(\boldsymbol{\xi}^{(k^*+1)}, \mathcal{G}^{(k^*+1)})$$

Consequently, the estimators could be obtained by $Q(\boldsymbol{\xi}^{(k^*)}, \mathcal{G}^{(k^*)})$.

(2) Local Minimality.

Next, we prove that the solution $(\boldsymbol{\xi}^{(k^*)}, \mathcal{G}^{(k^*)})$ is the local minimum of the function $\boldsymbol{\xi} \stackrel{\text{def}}{=} \operatorname{min}_{\mathcal{G}} Q(\boldsymbol{\xi}, \mathcal{G})$. Define $\tilde{\mathcal{G}} \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathcal{G}} Q(\boldsymbol{\xi}^{(k^*)} + \boldsymbol{\delta}, \mathcal{G})$, where $\boldsymbol{\delta}$ is a small perturbation. When $\|\boldsymbol{\delta}\|$ is small enough, we have that $\tilde{\mathcal{G}} = \mathcal{G}^{(k^*)}$ due to the group memberships' discreteness. Then, we have

$$\begin{aligned}
 \min_{\mathcal{G}} Q(\boldsymbol{\xi}^{(k^*)} + \boldsymbol{\delta}, \mathcal{G}) &= Q(\boldsymbol{\xi}^{(k^*)} + \boldsymbol{\delta}, \tilde{\mathcal{G}}) \\
 &= Q(\boldsymbol{\xi}^{(k^*)} + \boldsymbol{\delta}, \mathcal{G}^{(k^*)}) \geq Q(\boldsymbol{\xi}^{(k^*)}, \mathcal{G}^{(k^*)}) = \min_{\mathcal{G}} Q(\boldsymbol{\xi}^{(k^*)}, \mathcal{G}).
 \end{aligned}$$

This completes the proof.

Appendix F. Initialization

We provide the initialization procedure for the group memberships \mathcal{G}_l in the Algorithm A.3.

Appendix G. Additional Discussion of Mixed GTNAR Model (24)
G.1 Estimation Procedure

The element-wise model form can be given as

$$Y_{ij,t} = \gamma_{g_i^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} \sum_k \sum_m \frac{a_{im}^{(1)}}{n_{1i}} Y_{mk,t-1} \frac{a_{kj}^{(2)}}{n_{2j}}$$

Algorithm A.3 Initialization of the GTNAR Model

- 1: **Input:** $\{\mathcal{Y}_t, \mathbf{X}_t^{(l)}, \mathbf{W}^{(l)}, G_l\}$ for $1 \leq l \leq q$.
 - 2: Treat each node as a group for all dimensions $l \in [q]$, estimate $\hat{\boldsymbol{\theta}}^{(l)}$ and $\hat{\boldsymbol{\alpha}}$ by (13). Here $\hat{\boldsymbol{\theta}}^{(l)} \in \mathbb{R}^{N_l(p_l+1)}$ and $\hat{\boldsymbol{\alpha}} \in \mathbb{R}^{N_1 \times \dots \times N_q}$.
 - 3: Run the k -means clustering for the above estimators for $T_{\text{init}} = 3$ trials. For each trial $t = 1, \dots, T_{\text{init}}$, try the following two clustering types:
 - Type 1. (Clustering by time effect)**
 STEP 1. Clustering the nodes in the l th dimension by the mode- l matricization of *self-driven time effect* $\hat{\boldsymbol{\alpha}}_{(l)} \in \mathbb{R}^{N_l \times (\prod_{m \neq l} N_m)}$. Then one obtain the $\mathcal{G}_l^{1[t]}$ for $1 \leq l \leq q$.
 STEP 2. Calculate the loss in the t th trial for type 1 by $Q^{[t]}(\hat{\boldsymbol{\xi}}, \mathcal{G}^{1[t]})$, where $\mathcal{G}^{1[t]} = \{\mathcal{G}_l^{1[t]} : l \in [q]\}$ is the first type initial group memberships.
 - Type 2. (Clustering by network and covariate effect)**
 STEP 1. Clustering by *network effects* and *covariate effect* on $\hat{\boldsymbol{\theta}}^{(l)} = (\hat{\boldsymbol{\theta}}_{g^{(l)}}^{(l)} : g^{(l)} \in [G_l])$. Then one obtain the $\mathcal{G}_l^{2[t]}$ for $1 \leq l \leq q$.
 STEP 2. Calculate the loss in the t th trial for type 2 by $Q^{[t]}(\hat{\boldsymbol{\xi}}, \mathcal{G}^{2[t]})$, where $\mathcal{G}^{2[t]} = \{\mathcal{G}_l^{2[t]} : l \in [q]\}$ is the second type initial group memberships.
 - 4: Select the best initial trial $t^* = \operatorname{argmin}_t [\min\{Q^{[t]}(\hat{\boldsymbol{\xi}}, \mathcal{G}^{1[t]}), Q^{[t]}(\hat{\boldsymbol{\xi}}, \mathcal{G}^{2[t]})\}]$, and the corresponding initial memberships are denoted as $\mathcal{G}^{[0]}$.
 - 5: **Output:** Best initial memberships: $\mathcal{G}^{[0]}$.
-

$$\begin{aligned}
 & + \left(\lambda_{g_i^{(1)}}^{(1)} \sum_k \frac{a_{ik}^{(l)}}{n_{1i}} Y_{kj,t-1} + \lambda_{g_j^{(2)}}^{(2)} \sum_k \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,t-1} \right) \\
 & + \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)} + \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} \right) + \varepsilon_{ij,t}.
 \end{aligned}$$

For simplicity, denote the interactive parameters as $\boldsymbol{\psi} = (\boldsymbol{\psi}^{(1)\top}, \boldsymbol{\psi}^{(2)\top})^\top \in \mathbb{R}^{G_1+G_2}$ with $\boldsymbol{\psi}^{(1)} = (\gamma_1^{(1)}, \dots, \gamma_{G_1}^{(1)})^\top \in \mathbb{R}^{G_1}$ and $\boldsymbol{\psi}^{(2)} = (\gamma_1^{(2)}, \dots, \gamma_{G_2}^{(2)})^\top \in \mathbb{R}^{G_2}$. Recall that $\boldsymbol{\xi} = (\boldsymbol{\theta}^{(1)\top}, \boldsymbol{\theta}^{(2)\top}, \operatorname{vec}(\boldsymbol{\alpha})^\top)^\top \in \mathbb{R}^{G_1(p_1+1)+G_2(p_2+1)+G_1G_2}$ includes the other group parameters except for the interactive parameters $\boldsymbol{\psi}$. To estimate the unknown parameters, we minimize the following objective function,

$$\begin{aligned}
 Q(\boldsymbol{\xi}, \boldsymbol{\psi}, \mathcal{G}) = & \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} \sum_{t=1}^T \left\{ Y_{ij,t} - \gamma_{g_i^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} \sum_k \sum_m \frac{a_{im}^{(1)}}{n_{1i}} Y_{mk,t-1} \frac{a_{kj}^{(2)}}{n_{2j}} \right. \\
 & - \left(\lambda_{g_i^{(1)}}^{(1)} \sum_k \frac{a_{ik}^{(l)}}{n_{1i}} Y_{kj,t-1} + \lambda_{g_j^{(2)}}^{(2)} \sum_k \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,t-1} \right) \\
 & \left. - \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} \right) \right\}^2. \quad (\text{A.26})
 \end{aligned}$$

We utilize the following iterative algorithm.

- (1) **Update the group parameters $\boldsymbol{\xi}$ and $\boldsymbol{\psi}$.**

First, fix the group memberships $g_i^{(1)}$ and $g_j^{(2)}$ and then we update the group parameters. In this parameter updating step, we iteratively update $\boldsymbol{\psi}$ and $\boldsymbol{\xi}$ as follows, which can be conducted in a fast speed. To be more specific, given interactive network effects $\boldsymbol{\psi}$, we could calculate

$$\tilde{\mathbf{Y}}_t = \mathbf{Y}_t - (\boldsymbol{\Gamma}^{(1)} \mathbf{W}^{(1)}) \mathbf{Y}_{t-1} (\mathbf{W}^{(2)} \boldsymbol{\Gamma}^{(2)}).$$

Then the parameter $\boldsymbol{\xi}$ can be obtained by $\hat{\boldsymbol{\xi}} = (\mathbf{M})^{-1} \mathbf{b}$, whose specific expressions can be found in (A.7), with $\mathbb{Y}_{g^{(1)}, g^{(2)}, t}$ replaced by $\text{vec}(\tilde{\mathbf{Y}}_t^{\mathcal{R}_{g^{(1)}, \cdot}, \mathcal{R}_{g^{(2)}, \cdot}})$.

Subsequently, given $\boldsymbol{\xi}$, denote $\check{\mathbf{Y}}_t = \mathbf{Y}_t - \mathbf{L}^{(1)} \mathbf{W}^{(1)} \mathbf{Y}_{t-1} - \mathbf{Y}_{t-1} \mathbf{W}^{(2)} \mathbf{L}^{(2)} - \mathbf{A} \circ \mathbf{Y}_{t-1} - \boldsymbol{\beta}_{X_1, t}^{(1)} \mathbf{1}_{N_2}^\top - \mathbf{1}_{N_1} \boldsymbol{\beta}_{X_2, t}^{(2)\top}$. We need to minimize the following least squares problem,

$$\min_{\boldsymbol{\Gamma}^{(1)}, \boldsymbol{\Gamma}^{(2)}} \sum_t \|\check{\mathbf{Y}}_t - \boldsymbol{\Gamma}^{(1)} \mathbf{W}^{(1)} \mathbf{Y}_{t-1} \mathbf{W}^{(2)} \boldsymbol{\Gamma}^{(2)}\|_F^2.$$

Denote that $\tilde{\mathbf{X}}_t = \mathbf{W}^{(1)} \mathbf{Y}_{t-1} \mathbf{W}^{(2)} \boldsymbol{\Gamma}^{(2)} \in \mathbb{R}^{N_1 \times N_2}$. By the first order condition, we have that

$$\hat{\gamma}_{g^{(1)}}^{(1)} = \left(\sum_t \sum_{i \in \mathcal{R}_{g^{(1)}}} \sum_j \tilde{\mathbf{X}}_{ij, t} \check{\mathbf{Y}}_{ij, t} \right) / \sum_t \left\| \tilde{\mathbf{X}}_t^{\mathcal{R}_{g^{(1)}, \cdot}, \cdot} \right\|_F^2, \quad g^{(1)} \in [G_1] \quad (\text{A.27})$$

$$\hat{\gamma}_{g^{(2)}}^{(2)} = \left(\sum_t \sum_i \sum_{j \in \mathcal{R}_{g^{(2)}}} \tilde{\mathbf{X}}_{ij, t} \check{\mathbf{Y}}_{ij, t} \right) / \sum_t \left\| \tilde{\mathbf{X}}_t^{\cdot, \mathcal{R}_{g^{(2)}, \cdot}} \right\|_F^2, \quad g^{(2)} \in [G_2], \quad (\text{A.28})$$

where $\mathcal{R}_{g^{(1)}}^{(1)} = \{i \in [N_1] : g_i^{(1)} = g^{(1)}\}$ and $\mathcal{R}_{g^{(2)}}^{(2)} = \{j \in [N_2] : g_j^{(2)} = g^{(2)}\}$. Hence, we can iteratively update $\boldsymbol{\Gamma}^{(1)}$ and $\boldsymbol{\Gamma}^{(2)}$ using the above equations (A.27) and (A.28).

(2) Update the group memberships \mathcal{G}_1 and \mathcal{G}_2 .

Second, we consider given $\boldsymbol{\xi}$ and $\boldsymbol{\psi}$, updating the group memberships $\mathcal{G}_1 = (g_i^{(1)} : i \in [N_1])^\top$ and $\mathcal{G}_2 = (g_j^{(2)} : j \in [N_2])^\top$ iteratively. Given \mathcal{G}_2 , update

$$\begin{aligned} \hat{g}_i^{(1)} = \operatorname{argmin}_{g_i^{(1)} \in [G_1]} & \sum_{j=1}^{N_2} \sum_t \left\{ Y_{ij, t} - \gamma_{g_i^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} \sum_k \sum_m \frac{a_{im}^{(1)}}{n_{1i}} Y_{mk, t-1} \frac{a_{kj}^{(2)}}{n_{2j}} \right. \\ & - \left(\lambda_{g_i^{(1)}}^{(1)} \sum_k \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj, t-1} + \lambda_{g_j^{(2)}}^{(2)} \sum_k \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik, t-1} \right) \\ & \left. - \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij, (t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} \right) \right\}^2. \end{aligned} \quad (\text{A.29})$$

The parallel computation is implement to update \mathcal{G}_2 given \mathcal{G}_1 . To summarize, the iterative algorithm is shown in Algorithm A.4.

To estimate the group numbers for model (25), denoted as $\underline{G} = (G_1, G_2)$, we minimize the information criterion.

$$\text{QIC}(\underline{G}) = \log\{Q(\hat{\boldsymbol{\xi}}(\underline{G}), \hat{\boldsymbol{\psi}}(\underline{G}), \hat{\mathcal{G}}(\underline{G}))\} + \lambda(\underline{G}), \quad (\text{A.30})$$

where $\widehat{\boldsymbol{\xi}}_1(\underline{G})$, $\widehat{\boldsymbol{\psi}}(\underline{G})$ and $\widehat{\mathcal{G}}(\underline{G})$ are the estimated parameters when specifying the group numbers as G_1 and G_2 , $Q(\cdot, \cdot, \cdot)$ is the objective function defined in (A.26), and $\lambda(\underline{G})$ is the penalty function. Specifically, we set $\lambda(\underline{G}) = \kappa_1(G_1 + G_2)$ with κ_1 being a tuning parameter.

We summarize the estimation algorithm in the following Algorithm A.4.

Algorithm A.4 Estimation of the GTNAR Model with Interactive Term When $q = 2$

- 1: **Input:** $\{\mathbf{Y}_t, \mathbf{X}_t^{(l)}, \mathbf{W}^{(l)}, G_l\}$ for $1 \leq l \leq q$.
- 2: Obtain initial group memberships $\mathcal{G}_l^{[0]}$ according to Algorithm A.3. Let $\{\boldsymbol{\xi}^{[k]}, \mathcal{G}_l^{[k]}\}$ be the estimators and memberships in the k th iteration.
- 3: Repeat STEP 1 and STEP 2 for $k = 1, 2, \dots$ until convergence.

STEP 1. Given $\{\mathcal{G}_l^{[k-1]}\}$ for all layers $l \in [q]$

STEP (1A). Fix the interactive network effects $\boldsymbol{\Gamma}^{(1)[k-1]}$ and $\boldsymbol{\Gamma}^{(2)[k-1]}$, calculate $\widetilde{\mathbf{Y}}_t = \mathbf{Y}_t - (\boldsymbol{\Gamma}^{(1)[k-1]} \mathbf{W}^{(1)} \mathbf{Y}_{t-1} (\boldsymbol{\Gamma}^{(2)[k-1]}))$. Update $\boldsymbol{\xi}^{[k]} = (\mathbf{M}^{[k-1]})^{-1} \mathbf{b}^{[k-1]}$.

STEP (1B). Fix the estimated parameters $\boldsymbol{\xi}^{[k]}$. Given $\boldsymbol{\Gamma}^{(2)[k-1]}$, obtain $\boldsymbol{\Gamma}^{(1)[k]}$ using (A.27); then given $\boldsymbol{\Gamma}^{(1)[k]}$, obtain $\boldsymbol{\Gamma}^{(2)[k]}$ using (A.28).

STEP 2. Given $\gamma_{g_i^{(1)[k-1]}}^{(1)[k]}$, $\gamma_{g_j^{(2)[k-1]}}^{(2)[k]}$, $\lambda_{g_i^{(1)[k-1]}}^{(1)[k]}$, $\lambda_{g_j^{(2)[k-1]}}^{(2)[k]}$, $\alpha_{g_i^{(1)[k-1]} g_j^{(2)[k-1]}}^{[k]}$ and $\boldsymbol{\zeta}_{g_i^{(1)[k-1]}}^{(1)[k]}$, $\boldsymbol{\zeta}_{g_j^{(2)[k-1]}}^{(2)[k]}$, update the memberships $\mathcal{G}_1^{[k]}$ and $\mathcal{G}_2^{[k]}$ sequentially by,

$$\begin{aligned} g_i^{(1)[k]} &= \operatorname{argmin}_{g_i^{(1)} \in [G_1]} \sum_j \sum_t \left\{ Y_{ij,t} - \left(\gamma_{g_i^{(1)[k]}}^{(1)[k]} \gamma_{g_j^{(2)[k-1]}}^{(2)[k]} \sum_k \sum_m \frac{a_{im}^{(1)}}{n_{1i}} Y_{mk,t-1} \frac{a_{kj}^{(2)}}{n_{2j}} \right) \right. \\ &\quad - \left(\lambda_{g_i^{(1)[k]}}^{(1)[k]} \sum_k \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,t-1} + \lambda_{g_j^{(2)[k-1]}}^{(2)[k]} \sum_k \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,t-1} \right) \\ &\quad \left. - \alpha_{g_i^{(1)[k]} g_j^{(2)[k-1]}}^{[k]} Y_{ij,(t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)[k]}}^{(1)[k]} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)[k-1]}}^{(2)[k]} \right) \right\}^2 \end{aligned}$$

and

$$\begin{aligned} g_j^{(2)[k]} &= \operatorname{argmin}_{g_j^{(2)} \in [G_2]} \sum_i \sum_t \left\{ Y_{ij,t} - \left(\gamma_{g_i^{(1)[k]}}^{(1)[k]} \gamma_{g_j^{(2)[k]}}^{(2)[k]} \sum_k \sum_m \frac{a_{im}^{(1)}}{n_{1i}} Y_{mk,t-1} \frac{a_{kj}^{(2)}}{n_{2j}} \right) \right. \\ &\quad - \left(\lambda_{g_i^{(1)[k]}}^{(1)[k]} \sum_k \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,t-1} + \lambda_{g_j^{(2)[k]}}^{(2)[k]} \sum_k \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,t-1} \right) \\ &\quad \left. - \alpha_{g_i^{(1)[k]} g_j^{(2)[k]}}^{[k]} Y_{ij,(t-1)} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)[k]}}^{(1)[k]} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)[k]}}^{(2)[k]} \right) \right\}^2 \end{aligned}$$

- 4: **Output:** Final estimator and memberships: $\widehat{\boldsymbol{\xi}} = \boldsymbol{\xi}^{[K]}$, $\widehat{\boldsymbol{\psi}} = \boldsymbol{\psi}^{[K]}$, and $\widehat{\mathcal{G}} = \{\mathcal{G}_l^{[K]} : l = 1, 2\}$. Here K is the final number of iteration rounds.
-

G.2 Technical Assumptions

Denote

$$\mathcal{X}_{ij,t} = \left(\sum_{k=1}^{N_1} w_{ik}^{(1)} Y_{kj,t}, \mathbf{x}_{1t}^{(1)\top}, \sum_{k=1}^{N_2} w_{kj}^{(2)} Y_{ik,t}, \mathbf{x}_{jt}^{(2)\top}, Y_{ij,t-1}, \sum_k \sum_m w_{ik}^{(1)} Y_{km,t-1} w_{mj}^{(2)} \right) \in \mathbb{R}^{p_1+p_2+4}. \quad (\text{A.31})$$

By using $\mathcal{X}_{ij,t}$, we rewrite the objective function as $Q(\Theta) = \sum_{i,j,t} (Y_{ij,t} - \mathcal{X}_{ij,t}^\top \Theta_{ij})^2$. Denote $\widehat{\Theta}$ as the estimator of Θ by minimizing the loss function $Q(\Theta)$. We require the following assumptions.

Assumption A.11. (PARAMETER SPACE) *The parameter satisfies that $\|\Theta\|_{\max} < \infty$.*

Assumption A.12. (CONVEXITY) *Denote $\Sigma_{ij} = E(\mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top)$, where $\mathcal{X}_{ij,t}$ is defined in (A.31). Assume that $\tau_{\min}^1 \stackrel{\text{def}}{=} \min_{i,j} \Sigma_{ij} > 0$ is a constant.*

Assumption A.13. (STABILITY) *Suppose $\mathbb{Y}_0 \stackrel{\text{def}}{=} \text{vec}(\mathcal{Y}_0) = \mathbf{0}$, and assume that*

$$\max_{g^{(1)} \in [G_{1,0}], g^{(2)} \in [G_{2,0}]} \left| \lambda_{g^{(1)}}^{(1)0} + \lambda_{g^{(2)}}^{(2)0} + \alpha_{g^{(1)}g^{(2)}}^0 + \gamma_{g^{(1)}}^{(1)0} \gamma_{g^{(2)}}^{(2)0} \right| \leq \kappa_{\max} < 1,$$

where $G_{l,0}$ is the true number of groups in the l th dimension, and κ_{\max} is a positive constant.

Assumption A.14. (GROUP DIFFERENCE) *For the first dimension, assume*

$$\min_{g_1^{(1)} \neq g_2^{(1)}} \left[\left\| \boldsymbol{\theta}_{g_1^{(1)}}^{(1)0} - \boldsymbol{\theta}_{g_2^{(1)}}^{(1)0} \right\|^2 + \max_{g^{(2)} \in [G_2^0]} \left\{ \left| \alpha_{g_1^{(1)}g^{(2)}}^0 - \alpha_{g_2^{(1)}g^{(2)}}^0 \right|^2 + \left| \gamma_{g_1^{(1)}}^{(1)0} \gamma_{g^{(2)}}^{(2)0} - \gamma_{g_2^{(1)}}^{(1)0} \gamma_{g^{(2)}}^{(2)0} \right|^2 \right\} \right] \geq c_{\text{gap}},$$

where $c_{\text{gap}} \gg T^{-1}(\log(N_1 N_2))^2 (G_1 G_2)$ as $T \rightarrow \infty$. Here G_1, G_2 are finite. Assume the same holds for the group $g_1^{(2)} \neq g_2^{(2)}$ in the second dimension.

G.3 Numerical Studies

The estimation results in Table A.6 demonstrate excellent finite sample performance: all RMSEs decrease with increasing time length T or increasing sizes N_1, N_2 (e.g., RMSE of $\widehat{\boldsymbol{\lambda}}^{(1)}$ declines from 0.0131 at $T = 40$ to 0.0028 at $T = 80$), while the coverage probabilities in parentheses converge closely to the nominal 95% level at larger T . Furthermore, the group membership errors $\widehat{\eta}_1$ and $\widehat{\eta}_2$ remain near-zero across all configuration with consistently diminishing magnitudes at larger sample sizes, indicating exceptionally accurate group classification even for moderate sample sizes.

G.4 Technical Proofs

G.4.1 PROOF OF THEOREM 8 (I)

By using the $\mathcal{X}_{ij,t}$ defined in (A.31), we define the objective function,

$$Q(\Theta) = \sum_{i=1}^{N_1} \sum_{j=2}^{N_2} \sum_{t=1}^T (Y_{ij,t} - \mathcal{X}_{ij,t}^\top \Theta_{ij})^2 = \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} \sum_{t=1}^T Q_{ij}(\Theta_{ij}).$$

Table A.6: RMSEs of estimated parameters under when $G_{1,0} = G_{2,0} = 3$ with 100 replications. The performances are evaluated for different sample sizes N_1, N_2 and the time length T . The corresponding CPs are shown in the parenthesis.

G_1	G_2	N_1	N_2	T	$\widehat{\lambda}^{(1)}$	$\widehat{\lambda}^{(2)}$	$\widehat{\zeta}^{(1)}$	$\widehat{\zeta}^{(2)}$	$\widehat{\alpha}$	$\widehat{\Gamma}_{\text{norm}}$	$\widehat{\eta}_1$	$\widehat{\eta}_2$		
3	3	100	80	20	0.0131 (0.90)	0.0119 (0.93)	0.0137 (0.93)	0.0150 (0.92)	0.0369 (0.92)	3.0282	0.0122	0.0086		
				40	0.0079 (0.93)	0.0081 (0.95)	0.0089 (0.95)	0.0094 (0.95)	0.0131 (0.95)	1.2326	0.0005	0.0001		
				80	0.0055 (0.93)	0.0052 (0.95)	0.0064 (0.95)	0.0067 (0.95)	0.0091 (0.95)	0.5685	0.0000	0.0001		
				200	150	20	0.0061 (0.94)	0.0053 (0.93)	0.0070 (0.94)	0.0069 (0.94)	0.0095 (0.94)	2.5897	0.0004	0.0002
						40	0.0040 (0.95)	0.0039 (0.93)	0.0047 (0.95)	0.0048 (0.96)	0.0068 (0.95)	1.0622	0.0000	0.0002
						80	0.0027 (0.95)	0.0028 (0.93)	0.0032 (0.96)	0.0035 (0.96)	0.0047 (0.94)	0.6653	0.0000	0.0001

Given this objective function, we can follow the proof of Theorem 2 using the partition of the new parameter space. To obtain the conclusion, we use the Lemma 15 and Lemma 13 instead of the usage of Lemma 23 and Lemma 19 in the original proof of Theorem 2. Since the rest of the proof remains the same, we omit the details here.

G.4.2 PROOF OF THEOREM 8 (III)

We first declare some notations in the proof. Denote $\widetilde{\gamma}_{g^{(1)}}^{(1)} = (\gamma_{g^{(1)}}^{(1)} \gamma_{g^{(2)}}^{(2)} : g^{(2)} \in [G_2])^\top \in \mathbb{R}^{G_2}$ and $\widetilde{\gamma}_{g^{(2)}}^{(2)} = (\gamma_{g^{(1)}}^{(2)} \gamma_{g^{(2)}}^{(2)} : g^{(1)} \in [G_1])^\top \in \mathbb{R}^{G_1}$. Denote $\widetilde{\gamma} = (\widetilde{\gamma}_{g^{(1)}}^{(1)\top} : g^{(1)} \in [G_1])^\top \in \mathbb{R}^{G_1 G_2}$. Further denote that $\widetilde{\alpha} = (\text{vec}(\alpha_{g^{(1)}})^\top : g^{(1)} \in [G_1])^\top \in \mathbb{R}^{G_1 G_2}$. Recall that $\theta^{(1)} = (\theta_{g^{(1)}}^{(1)\top} : g^{(1)} \in [G_1])^\top$. Denote $\xi_{g^{(1)}}^{(1)} = (\theta_{g^{(1)}}^{(1)\top}, \text{vec}(\alpha_{g^{(1)}})^\top, \widetilde{\gamma}_{g^{(1)}}^{(1)\top})^\top$ as the parameters of the $g^{(1)}$ th group ($g^{(1)} \in [G_1]$). Correspondingly, denote $\xi_{g^{(1)0}}^{(1)0}$ as the true parameters of the true group $g^{(1)0} \in [G_{1,0}]$. Let $\xi^{(1)} = (\theta^{(1)\top}, \widetilde{\alpha}^\top, \widetilde{\gamma}^\top)^\top$. In parallel, we can define $\xi_{g^{(2)}}^{(2)} = (\theta_{g^{(2)}}^{(2)\top}, \widetilde{\alpha}^\top, \widetilde{\gamma}^\top)^\top$ and $\xi^{(2)} = (\theta^{(2)\top}, \widetilde{\alpha}^\top, \widetilde{\gamma}^\top)^\top$. At last, denote $\xi = (\theta^{(1)\top}, \theta^{(2)\top}, \widetilde{\alpha}^\top, \widetilde{\gamma}^\top)^\top$. In the following, we prove the theorem for the first dimension, and the second dimension group membership consistency can be proved similarly. For $g^{(1)} \in [G_1]$ and $g^{(1)0} \in [G_{1,0}]$, define

$$\begin{aligned} \mathcal{L}^{(1)}(\xi_{g^{(1)}}^{(1)}, \xi_{g^{(1)0}}^{(1)0}; \mathcal{G}_2, \mathcal{G}_2^0) &= \|\theta_{g^{(1)}}^{(1)} - \theta_{g^{(1)0}}^{(1)0}\|^2 \\ &+ \frac{1}{N_2} \sum_j \left\{ |\alpha_{g^{(1)}g_j^{(2)}} - \alpha_{g^{(1)0}g_j^{(2)0}}|^2 + |\gamma_{g^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} - \gamma_{g^{(1)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0}|^2 \right\}. \end{aligned} \quad (\text{A.32})$$

For notation simplicity, denote $\mathcal{L}_{g^{(1)}, g^{(1)0}}^{(1)} \stackrel{\text{def}}{=} \mathcal{L}^{(1)}(\xi_{g^{(1)}}^{(1)}, \xi_{g^{(1)0}}^{(1)0}; \mathcal{G}_2, \mathcal{G}_2^0)$. Define the distance

$$d_L^{(1)}(\xi^{(1)}, \xi^{(1)0}; \mathcal{G}_2, \mathcal{G}_2^0) = \max \left\{ \max_{g^{(1)0} \in [G_{1,0}]} \min_{g^{(1)} \in [G_1]} \mathcal{L}_{g^{(1)}, g^{(1)0}}^{(1)}, \max_{g^{(1)} \in [G_1]} \min_{g^{(1)0} \in [G_{1,0}]} \mathcal{L}_{g^{(1)}, g^{(1)0}}^{(1)} \right\}.$$

In addition, given parameters $\boldsymbol{\xi}$, denote the estimated group memberships in the second dimension as $\widehat{\mathcal{G}}_2(\boldsymbol{\xi}) = (\widehat{g}_j^{(2)}(\boldsymbol{\xi}) : j \in [N_2])^\top$, define

$$\mathcal{N}_\eta^{(1)} = \{\boldsymbol{\xi} : d_L^{(1)}(\boldsymbol{\xi}^{(1)}, \boldsymbol{\xi}^{(1)0}; \widehat{\mathcal{G}}_2(\boldsymbol{\xi}), \mathcal{G}_2^0) < \eta\}$$

and denote

$$\begin{aligned} \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g^{(1)0}, \mathcal{G}_2^0) &= \left\{ g^{(1)} \in [G_1] : \right. \\ &\left. \|\boldsymbol{\theta}_{g^{(1)}}^{(1)} - \boldsymbol{\theta}_{g^{(1)0}}^{(1)0}\|^2 + \frac{1}{N_2} \sum_j \left\{ |\alpha_{g^{(1)}\widehat{g}_j^{(2)}(\boldsymbol{\xi})} - \alpha_{g^{(1)0}g_j^{(2)0}}^0|^2 + |\gamma_{g^{(1)}\widehat{g}_j^{(2)}(\boldsymbol{\xi})}^{(1)} - \gamma_{g^{(1)0}g_j^{(2)0}}^{(1)0}|^2 \right\} \leq \eta \right\}. \end{aligned}$$

Based on the above notations, we first give the following lemma, whose proof can be found in Appendix G.4.3.

Lemma 12. *Under Assumptions 3, 4, 7, A.11–A.14, and assume that $G_1 \geq G_{1,0}$ and $G_2 \geq G_{2,0}$. For Θ satisfying the condition*

$$d(\Theta, \Theta^0) = O_p\{T^{-1}(\log(N_1N_2))^2\} = o_p(c_{\text{gap}}) \quad (\text{A.33})$$

as $N_l \rightarrow \infty$, $l = 1, 2$, the following conclusions hold,

(i) For all $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(1)}$ with $\eta < c_{\text{gap}}c_\pi/4$, then we have that $\{\mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g^{(1)0}, \mathcal{G}_2^0) : g^{(1)0} \in [G_{1,0}]\}$ is a partition of $[G_1]$.

(ii) For all $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(1)}$ with $\eta \leq \tau_{\min}^1 c_{\text{gap}} c_\pi / (8(\tau_{\min}^1 + \tau_{\max}^1))$, define the event $\mathcal{O} = \{\widehat{g}_i^{(1)}(\boldsymbol{\xi}) \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0) : i \in [N_1]\}$. Then we have $P(\mathcal{O}^c) \leq C \exp(-c_1 \sqrt{T} c_{\text{gap}} + c_2 m + \log N_1 + \log N_2)$, where C, c_1, c_2 are positive constants, and $m = p_1 + p_2 + 4$.

Next, by using Lemma 12, we first show that $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(1)}$ when $\eta = \tau_{\min}^1 c_{\text{gap}} c_\pi / (8(\tau_{\min}^1 + \tau_{\max}^1))$ with probability tending to 1, then we use (i) and (ii) in Lemma 12 to obtain the final results.

Step 1. $P(\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(1)}) \rightarrow 1$.

Using the similar procedure in the Step 1 of proof of Theorem 5, we can show that $d_L^{(1)}(\widehat{\boldsymbol{\xi}}^{(1)}, \boldsymbol{\xi}^{(1)0}; \widehat{\mathcal{G}}_2(\widehat{\boldsymbol{\xi}}), \mathcal{G}_2^0) = O_p\{T^{-1}(\log(N_1N_2))^2\}$. By the condition that

$$c_{\text{gap}} \gg T^{-1}(\log(N_1N_2))^2,$$

we have $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(1)}$ with probability tending to 1.

Step 2. Conclusion of Theorem 8 (iii).

For $\widetilde{g}^{(1)} \in [G_1]$, suppose there are $i_1, i_2 \in \widehat{\mathcal{R}}_{\widetilde{g}^{(1)}}^{(1)}$, then $\widehat{g}_{i_1}^{(1)} = \widehat{g}_{i_2}^{(2)} = \widetilde{g}^{(1)}$. By Lemma 12(ii), we have $\widetilde{g}^{(1)} \in \mathcal{A}_\eta^{(1)}(\widehat{\boldsymbol{\xi}}, g_{i_1}^{(1)0}, \mathcal{G}_2^0)$ and $\widetilde{g}^{(1)} \in \mathcal{A}_\eta^{(1)}(\widehat{\boldsymbol{\xi}}, g_{i_2}^{(1)0}, \mathcal{G}_2^0)$ hold with probability tending to 1. By Lemma 12(i), we know that $\mathcal{A}_\eta^{(1)}(\widehat{\boldsymbol{\xi}}, \cdot, \mathcal{G}_2^0)$ is a partition of $[G_1]$, hence $g_{i_1}^{(1)0} = g_{i_2}^{(1)0}$. Then, by setting $g^{(1)} = g_{i_1}^{(1)0} = g_{i_2}^{(1)0}$, we have $i_1, i_2 \in \mathcal{R}_{g^{(1)}}^{(1)}$.

G.4.3 PROOF OF LEMMA 12

(1) Proof of (i).

For all $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(1)}$, by the definition of $\mathcal{N}_\eta^{(1)}$ and the group set definition $\mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g^{(1)0}, \mathcal{G}_2^0)$, we have $\cup_{g^{(1)0}=1}^{G_{1,0}} \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g^{(1)0}, \mathcal{G}_2^0) = [G_1]$. Then we prove $\mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_1^{(1)0}, \mathcal{G}_2^0) \cap \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_2^{(1)0}, \mathcal{G}_2^0) = \emptyset$ for $g_1^{(1)0}, g_2^{(1)0} \in [G_{1,0}]$ and $g_1^{(1)0} \neq g_2^{(1)0}$. We use the contradiction. Assume there exists $g_3^{(1)} \in [G_1]$ s.t. $g_3^{(1)} \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_1^{(1)0}, \mathcal{G}_2^0) \cap \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_2^{(1)0}, \mathcal{G}_2^0)$. Then under Assumption A.14, we have

$$\begin{aligned}
 c_{\text{gap}} &\leq \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_2^{(1)0}}^{(1)0} \right\|^2 + \max_{g^{(2)0} \in [G_2^0]} \left\{ \left| \alpha_{g_1^{(1)0} g^{(2)0}}^0 - \alpha_{g_2^{(1)0} g^{(2)0}}^0 \right|^2 + \left| \gamma_{g_1^{(1)0} g^{(2)0}}^{(1)0} \gamma_{g^{(2)0}}^{(2)0} - \gamma_{g_2^{(1)0} g^{(2)0}}^{(1)0} \gamma_{g^{(2)0}}^{(2)0} \right|^2 \right\} \\
 &\leq 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 + 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 \\
 &\quad + \max_{g^{(2)0} \in [G_{2,0}]} \left\{ \frac{1}{\pi_{g^{(2)0}, N_2}^{(2)}} \sum_j \left| \alpha_{g_1^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_2^{(1)0} g_j^{(2)0}}^0 \right|^2 I(g_j^{(2)0} = g^{(2)0}) \right\} \\
 &\quad + \max_{g^{(2)0} \in [G_{2,0}]} \left\{ \frac{1}{\pi_{g^{(2)0}, N_2}^{(2)}} \sum_j \left| \gamma_{g_1^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} - \gamma_{g_2^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} \right|^2 I(g_j^{(2)0} = g^{(2)0}) \right\} \\
 &\leq 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 + 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 \\
 &\quad + \max_{g^{(2)0} \in [G_{2,0}]} \left(\frac{1}{\pi_{g^{(2)0}, N_2}^{(2)}} \right) \sum_j \left\{ \left| \alpha_{g_1^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_2^{(1)0} g_j^{(2)0}}^0 \right|^2 + \left| \gamma_{g_1^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} - \gamma_{g_2^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} \right|^2 \right\} \\
 &\leq 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 + 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 \\
 &\quad + \left(\frac{1}{\min_{g^{(2)0} \in [G_{2,0}]} \pi_{g^{(2)0}, N_2}^{(2)}} \right) \sum_j \left\{ \left| \alpha_{g_1^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_2^{(1)0} g_j^{(2)0}}^0 \right|^2 \right. \\
 &\quad \quad \left. + \left| \gamma_{g_1^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} - \gamma_{g_2^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} \right|^2 \right\} \\
 &\leq \left(\frac{2}{\min_{g^{(2)0} \in [G_{2,0}]} \pi_{g^{(2)0}, N_2}^{(2)}} \right) \left\{ \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 + 2 \left\| \boldsymbol{\theta}_{g_1^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_3^{(1)}}^{(1)0} \right\|^2 \right. \\
 &\quad + \frac{1}{N_2} \sum_j \left| \alpha_{g_1^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_3^{(1)} \hat{g}_j^{(2)}(\boldsymbol{\xi})}^0 \right|^2 + \frac{1}{N_2} \sum_j \left| \alpha_{g_2^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_3^{(1)} \hat{g}_j^{(2)}(\boldsymbol{\xi})}^0 \right|^2 \\
 &\quad \left. + \frac{1}{N_2} \sum_j \left| \gamma_{g_1^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} - \gamma_{g_3^{(1)} \hat{g}_j^{(2)}(\boldsymbol{\xi})}^{(1)0} \gamma_{\hat{g}_j^{(2)}(\boldsymbol{\xi})}^{(2)0} \right|^2 + \frac{1}{N_2} \sum_j \left| \gamma_{g_2^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0} - \gamma_{g_3^{(1)} \hat{g}_j^{(2)}(\boldsymbol{\xi})}^{(1)0} \gamma_{\hat{g}_j^{(2)}(\boldsymbol{\xi})}^{(2)0} \right|^2 \right\} \\
 &\leq \frac{4\eta}{\min_{g^{(2)0} \in [G_{2,0}]} \pi_{g^{(2)0}, N_2}^{(2)}},
 \end{aligned}$$

where the last line holds due to $g_3^{(1)} \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_1^{(1)0}, \mathcal{G}_2^0) \cap \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_2^{(1)0}, \mathcal{G}_2^0)$. However, this conclusion contradicts $\eta < c_{\text{gap}} c_\pi / 4$ as $N_2 \rightarrow \infty$. Therefore, we prove that $\mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}; g_1^{(1)0}, \mathcal{G}_2^0) \cap \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}; g_2^{(1)0}, \mathcal{G}_2^0) = \emptyset$, which completes the proof of (i).

(2) Proof of (ii).

Denote $\boldsymbol{\xi}_{g_i^{(1)}}^{(1)} = (\boldsymbol{\theta}_{g_i^{(1)}}^{(1)\top}, \text{vec}(\boldsymbol{\alpha}_{g_i^{(1)}})^{\top}, \tilde{\boldsymbol{\gamma}}_{g_i^{(1)}}^{(1)\top})$. Define the loss function

$$\begin{aligned} Q_i(\boldsymbol{\xi}_{g_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2) &= \sum_j \sum_t \left\{ Y_{ij,t} - \gamma_{g_i^{(1)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} \sum_k \sum_m w_{im}^{(1)} Y_{mk,t-1} w_{kj}^{(2)} \right. \\ &\quad - \lambda_{g_i^{(1)}}^{(1)} \sum_k w_{ik}^{(1)} Y_{kj,t-1} - \lambda_{g_j^{(2)}}^{(2)} \sum_k w_{kj}^{(2)} Y_{ik,t-1} \\ &\quad \left. - \alpha_{g_i^{(1)} g_j^{(2)}}^{(1,2)} Y_{ij,t-1} - \left(\mathbf{x}_{it}^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + \mathbf{x}_{jt}^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} \right) \right\}^2. \end{aligned} \quad (\text{A.34})$$

For simplicity, for any given $\boldsymbol{\xi}$ satisfying $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(1)}$, denote $\widehat{g}_i^{(1)}(\boldsymbol{\xi})$ as $\widehat{g}_i^{(1)}$. Hence we have that for any $\widetilde{g}^{(1)} \neq g^{(1)}$,

$$\{\widehat{g}_i^{(1)} = g^{(1)}\} \subseteq \{Q_i(\boldsymbol{\xi}_{g^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi})) \leq Q_i(\boldsymbol{\xi}_{\widetilde{g}^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi}))\}.$$

Hence, for $\widetilde{g}_i^{(1)} \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)$,

$$\begin{aligned} I(\widehat{g}_i^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)) &= \sum_{g^{(1)=1}^{G_1}} I(g^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)) I(\widehat{g}_i^{(1)} = g^{(1)}) \\ &\leq \sum_{g^{(1)=1}^{G_1}} I(g^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)) I\{Q_i(\boldsymbol{\xi}_{g^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi}))\} \stackrel{\text{def}}{=} \sum_{g^{(1)=1}^{G_1}} W_{ig^{(1)}}(\boldsymbol{\xi}). \end{aligned}$$

On one hand, for all $g^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)$ and $g^{(1)} \in [G_1]$, there exists a $g_{i_2}^{(1)0} \neq g_i^{(1)0}$, such that $g^{(1)} \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_{i_2}^{(1)0}, \mathcal{G}_2^0)$ due to the conclusion of (i). Then we have

$$\begin{aligned} &\|\boldsymbol{\theta}_{g^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_i^{(1)0}}^{(1)0}\|^2 + \frac{1}{N_2} \sum_j \left\{ |\alpha_{g^{(1)} \widehat{g}_j^{(2)}}(\boldsymbol{\xi}) - \alpha_{g_i^{(1)0} g_j^{(2)0}}^0|^2 + |\gamma_{g^{(1)} \widehat{g}_j^{(2)}}^{(1)}(\boldsymbol{\xi}) - \gamma_{g_i^{(1)0} g_j^{(2)0}}^{(1)0}(\boldsymbol{\xi})|^2 \right\} \\ &\geq \frac{1}{2} \|\boldsymbol{\theta}_{g_{i_2}^{(1)0}}^{(1)0} - \boldsymbol{\theta}_{g_i^{(1)0}}^{(1)0}\|^2 + \frac{1}{2N_2} \sum_j \left\{ |\alpha_{g_{i_2}^{(1)0} g_j^{(2)0}}^0 - \alpha_{g_i^{(1)0} g_j^{(2)0}}^0|^2 + |\gamma_{g_{i_2}^{(1)0} g_j^{(2)0}}^{(1)0} - \gamma_{g_i^{(1)0} g_j^{(2)0}}^{(1)0}|^2 \right\} \\ &\quad - \left[\|\boldsymbol{\theta}_{g^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_{i_2}^{(1)0}}^{(1)0}\|^2 + \frac{1}{N_2} \sum_j \left\{ |\alpha_{g^{(1)} \widehat{g}_j^{(2)}}(\boldsymbol{\xi}) - \alpha_{g_{i_2}^{(1)0} g_j^{(2)0}}^0|^2 + |\gamma_{g^{(1)} \widehat{g}_j^{(2)}}^{(1)}(\boldsymbol{\xi}) - \gamma_{g_{i_2}^{(1)0} g_j^{(2)0}}^{(1)0}(\boldsymbol{\xi})|^2 \right\} \right] \\ &\geq c_{\text{gap}} c_\pi / 2 - \eta, \end{aligned}$$

when $N_2 \rightarrow \infty$, where c_{gap} and c_π are defined in Assumption A.14 and Assumption 7, respectively. By Lemma 15, it holds for any $g^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)$,

$$\frac{1}{N_2 T} \left\{ Q_i^*(\boldsymbol{\xi}_{g^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi})) - Q_i^*(\boldsymbol{\xi}_{g_i^{(1)0}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2^0) \right\} \geq \tau_{\min}^1 (c_{\text{gap}} c_\pi / 2 - \eta). \quad (\text{A.35})$$

On the other hand, for $\widetilde{g}_i^{(1)} \in \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)$, it holds that

$$\frac{1}{N_2 T} \left\{ Q_i^*(\boldsymbol{\xi}_{\widetilde{g}_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi})) - Q_i^*(\boldsymbol{\xi}_{g_i^{(1)0}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2^0) \right\}$$

$$\begin{aligned}
 &\leq \tau_{\max}^1 \left\{ \|\boldsymbol{\theta}_{\widehat{g}_i^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_i^{(1)0}}^{(1)0}\|^2 + \frac{1}{N_2} \sum_j \|\boldsymbol{\theta}_{\widehat{g}_j^{(2)}(\boldsymbol{\xi})}^{(1)} - \boldsymbol{\theta}_{g_j^{(2)0}}^{(1)0}\|^2 \right. \\
 &\quad \left. + \frac{1}{N_2} \sum_j (|\alpha_{\widehat{g}_i^{(1)}\widehat{g}_j^{(2)}(\boldsymbol{\xi})} - \alpha_{g_i^{(1)0}g_j^{(2)0}}^0|^2 + |\gamma_{\widehat{g}_i^{(1)}\widehat{g}_j^{(2)}(\boldsymbol{\xi})}^{(1)}\gamma_{\widehat{g}_j^{(2)}(\boldsymbol{\xi})}^{(2)} - \gamma_{g_i^{(1)0}g_j^{(2)0}}^{(1)0}\gamma_{g_j^{(2)0}}^{(2)0}|^2) \right\} \\
 &\leq \tau_{\max}^1 (d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) + \eta) \leq \tau_{\max}^1 \{CT^{-1}(\log N_1 + \log N_2)^2 + \eta\}, \tag{A.36}
 \end{aligned}$$

with probability tending to 1, where the first inequality holds due to Lemma 15, and the last inequality holds due to (A.33). Together with (A.35), we have that with probability tending to 1,

$$\begin{aligned}
 &\frac{1}{N_2 T} Q_i^*(\boldsymbol{\xi}_{g^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi})) - Q_i^*(\boldsymbol{\xi}_{\widehat{g}_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \widehat{\mathcal{G}}_2(\boldsymbol{\xi})) \\
 &\geq \tau_{\min}^1 (c_{\text{gap}} c_\pi / 2 - \eta) - \tau_{\max}^1 \{CT^{-1}(\log N_1 + \log N_2)^2 + \eta\} \stackrel{\text{def}}{=} \epsilon_\eta.
 \end{aligned}$$

Then we apply the similar techniques as in (A.93) and (A.94), to obtain that

$$\begin{aligned}
 &P \left\{ \sup_i I(\widehat{g}_i^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)) \right\} \\
 &\leq \sum_{g^{(1)=1}^{G_1}} P \left\{ \sup_i W_{ig^{(1)}}(\boldsymbol{\xi}) > 0 \right\} \\
 &\leq G_1 \exp \left\{ -C_2 \min(T\epsilon_\eta^2, \sqrt{T}\epsilon_\eta) + C_3 m + \log N_1 + \log N_2 \right\}.
 \end{aligned}$$

by Lemma 13. Since $c_{\text{gap}} \gg T^{-1}(\log(N_1 N_2))^2$, by the condition on η , we have $\eta < \{\tau_{\min}^1 c_{\text{gap}} c_\pi / 4 - \tau_{\max}^1 CT^{-1}(\log(N_1 N_2))^2\} / (\tau_{\min}^1 + \tau_{\max}^1)$. Hence, by the definition of ϵ_η , we have $\epsilon_\eta > \tau_{\min}^1 c_{\text{gap}} c_\pi / 4$ as $N_2, T \rightarrow \infty$. This leads to that

$$P \left\{ \sup_i I(\widehat{g}_i^{(1)} \notin \mathcal{A}_\eta^{(1)}(\boldsymbol{\xi}, g_i^{(1)0}, \mathcal{G}_2^0)) \right\} \leq C \exp \left\{ -c_1 \sqrt{T} c_{\text{gap}} + c_2 m + \log N_1 + \log N_2 \right\},$$

where $m = p_1 + p_2 + 4$.

G.4.4 TECHNICAL LEMMAS FOR INTERACTIVE MODEL

Lemma 13. *Under Assumptions 3–4 and A.11–A.13, we have*

$$\begin{aligned}
 &P \left\{ \sup_{\|\boldsymbol{\Theta}_{ij}\|_{\max} \leq R} \left| T^{-1} Q_{ij}(\boldsymbol{\Theta}_{ij}) - Q_{ij}^*(\boldsymbol{\Theta}_{ij}) \right| \geq x \right\} \\
 &\leq C_1 \exp \left\{ -C_2 \min(Tx^2, \sqrt{T}x) + C_3 m \right\} \tag{A.37}
 \end{aligned}$$

Further recall that $S_{ij}(\boldsymbol{\Theta}_{ij}) = Q_{ij}(\boldsymbol{\Theta}_{ij}) - Q_{ij}(\boldsymbol{\Theta}_{ij}^0)$, we have that

$$P \left\{ \sup_{\boldsymbol{\Theta}_{ij}} T^{-1} \frac{|S_{ij}(\boldsymbol{\Theta}_{ij}) - S_{ij}^*(\boldsymbol{\Theta}_{ij})|}{d_{ij}(\boldsymbol{\Theta}_{ij}, \boldsymbol{\Theta}_{ij}^0)} > x \right\} \leq C_1 \exp \left\{ -C_2 \min(Tx^2, \sqrt{T}x) + C_3 m \right\}. \tag{A.38}$$

with $m = p_1 + p_2 + 4$, and C_1, C_2, C_3 are positive constants, and

$$\begin{aligned}
 d_{ij}(\boldsymbol{\Theta}_{ij}, \boldsymbol{\Theta}_{ij}^0) &= \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_i^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_i^{(1)}}^{(1)}\|^2 + \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_j^{(2)}}^{(2)} - \boldsymbol{\theta}_{g_j^{(2)}}^{(2)}\|^2 \\
 &\quad + |\widehat{\alpha}_{\widehat{g}_i^{(1)}\widehat{g}_j^{(2)}} - \alpha_{g_i^{(1)}g_j^{(2)}}|^2 + |\widehat{\gamma}_{\widehat{g}_i^{(1)}\widehat{g}_j^{(2)}}^{(1)}\widehat{\gamma}_{\widehat{g}_j^{(2)}}^{(2)} - \gamma_{g_i^{(1)}g_j^{(2)}}^{(1)}\gamma_{g_j^{(2)}}^{(2)}|^2.
 \end{aligned}$$

Proof
1. Proof of (A.37).

Note that

$$T^{-1}Q_{ij}(\Theta_{ij}) = T^{-1} \sum_t (\varepsilon_{ij,t} + \mathcal{X}_{ij,t}^\top \Theta_{ij}^0 - \mathcal{X}_{ij,t}^\top \Theta_{ij})^2,$$

where $\mathcal{X}_{ij,t}$ is defined in (A.31). We use the similar procedure in the proof of Lemma 17. The three core steps in the proof of Lemma 17 are (A.68)–(A.70). While the implementation of (A.68) remains the same, the interactive network terms in $\mathcal{X}_{ij,t}$ and Θ_{ij} cause changes in using (A.69) and (A.70). We take the change in (A.70) for example.

Specifically, the diagonal elements of $T^{-1} \sum_t \mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top$ take the forms $T^{-1} \sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w}$, $T^{-1} \sum_t \|\mathbf{x}_{it}^{(1)}\|^2$, and $T^{-1} \sum_t \|\mathbf{x}_{jt}^{(2)}\|^2$. For the additive network term, \mathbf{w} is replaced by $\mathbf{W}_{i \cdot}^{(1)\top} \otimes \mathbf{1}_{N_2} \in \mathbb{R}^{N_1 \times N_2}$, while for the interactive network term, \mathbf{w} is replaced by $\mathbf{W}_{i \cdot}^{(1)\top} \otimes \mathbf{W}_{\cdot j}^{(2)} \in \mathbb{R}^{N_1 \times N_2}$. Both types of network terms satisfy that $\|\mathbf{w}\|_1 = 1$. Hence, for $T^{-1} \sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w}$, we apply Lemma 18 to prove the concentration inequality. For $T^{-1} \sum_t \|\mathbf{x}_{it}^{(1)}\|^2$, and $T^{-1} \sum_t \|\mathbf{x}_{jt}^{(2)}\|^2$, we use Lemma 18 and Assumptions 4.

Then, similar with the Step 2 in the proof of Lemma 17, by taking union upper bound for all $\|\Theta_{ij}\|_{\max} \leq R$, we obtain the final conclusion.

2. Proof of (A.38).

Note that

$$S_{ij}(\Theta_{ij}) - S_{ij}^*(\Theta_{ij}) = [(\Theta_{ij} - \Theta_{ij}^0)^\top \mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top (\Theta_{ij} - \Theta_{ij}^0)] - E[(\Theta_{ij} - \Theta_{ij}^0)^\top \mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top (\Theta_{ij} - \Theta_{ij}^0)] \quad (\text{A.39})$$

$$+ 2\mathcal{X}_{ij,t}^\top (\Theta_{ij} - \Theta_{ij}^0) \varepsilon_{ij,t}. \quad (\text{A.40})$$

We borrow the idea to prove Lemma 17 to derive the results. For instance, we can apply (A.66) on (A.39) by using the fact that $\|\Theta_{ij} - \Theta_{ij}^0\| / \sqrt{d_{ij}(\Theta_{ij}, \Theta_{ij}^0)} \leq cR$ for some constant c . Similar techniques can be applied on proof of (A.40). \blacksquare

Lemma 14. *Under Assumptions 3–4 and A.11–A.13, for vector $\Delta \in \mathbb{R}^{p_1+p_2+4}$ in the parameter space of the interactive model (25), we have*

$$\begin{aligned} & P \left\{ \sup_i \sup_{\|\Delta\|^2 \leq \omega^2} (N_2 T)^{-1} \left| \sum_j \sum_t \left\{ \Delta^\top \mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top) \Delta] \right\} \right| \geq x \right\} \\ & \leq C_1 \exp \left\{ -C_2 \min \left(T x^2 / \omega^2, \sqrt{T} x / \omega \right) + C_3 m + \log(N_1 N_2) \right\}, \\ & P \left\{ \sup_i \sup_{\|\Delta\|^2 \leq \omega^2} (N_2 T)^{-1} \left| \sum_j \sum_t \mathcal{X}_{ij,t}^\top \Delta \varepsilon_{ij,t} \right| \geq x \right\} \\ & \leq C_1 \exp \left\{ -C_2 \min \left(T x^2 / \omega^2, \sqrt{T} x / \omega \right) + C_3 m + \log(N_1 N_2) \right\}, \end{aligned}$$

where $\mathcal{X}_{ij,t}$ is defined in equation (A.31), $m = p_1 + p_2 + 4$, C_1, C_2, C_3 are positive constants.

Proof The proof is similar with the proof of Lemma 20. The key difference lies in the expression of $\mathcal{X}_{ij,t}$. The technical details can refer to the proof of Lemma 13, which is omitted here. \blacksquare

Lemma 15. *Under Assumption A.12, it holds that*

$$\begin{aligned} \tau_{\min}^1 d_i(\Theta_{i,1}, \Theta_i^0) &\leq \frac{1}{N_2 T} \left\{ Q_i^*(\boldsymbol{\xi}_{g_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2) - Q_i^*(\boldsymbol{\xi}_{g_i^{(1)0}}^{(1)0}; \boldsymbol{\xi}_{g^{(2)0}}^{(2)0}, \mathcal{G}_2^0) \right\} \\ &\leq \tau_{\max}^1 d_i(\Theta_i, \Theta_i^0), \\ \tau_{\min}^1 d(\Theta, \Theta^0) &\leq \frac{1}{N_2 T} \left\{ Q^*(\Theta) - Q^*(\Theta^0) \right\} \leq \tau_{\max}^1 d(\Theta, \Theta^0) \end{aligned}$$

where τ_{\min}^1 is defined in Assumption A.12, and τ_{\max}^1 is defined in Lemma 16. Besides, $\Theta_i = (\Theta_{ij}^\top : j \in [N_2]) \in \mathbb{R}^{N_2 \times m}$, $\Theta_{ij} = (\boldsymbol{\theta}_{g_i^{(1)}}^{(1)\top}, \boldsymbol{\theta}_{g_j^{(2)}}^{(2)\top}, \alpha_{g_i^{(1)} g_j^{(2)}}, \gamma_{g_i^{(1)} g_j^{(2)}})^\top \in \mathbb{R}^m$, and

$$\begin{aligned} d_i(\Theta_i, \Theta_i^0) &= N_2^{-1} \sum_j \|\Theta_{ij} - \Theta_{ij}^0\|^2 \\ &= \|\boldsymbol{\theta}_{g_i^{(1)}}^{(1)} - \boldsymbol{\theta}_{g_i^{(1)0}}^{(1)0}\|^2 + N_2^{-1} \sum_j \|\boldsymbol{\theta}_{g_j^{(2)}}^{(2)} - \boldsymbol{\theta}_{g_j^{(2)0}}^{(2)0}\|^2 \\ &\quad + N_2^{-1} \sum_j \left\{ |\alpha_{g_i^{(1)} g_j^{(2)}} - \alpha_{g_i^{(1)0} g_j^{(2)0}}^0|^2 + |\gamma_{g_i^{(1)} g_j^{(2)}}^{(1)} \gamma_{g_j^{(2)}}^{(2)} - \gamma_{g_i^{(1)0} g_j^{(2)0}}^{(1)0} \gamma_{g_j^{(2)0}}^{(2)0}|^2 \right\}. \end{aligned}$$

Proof By the definition of loss function $Q_i(\boldsymbol{\xi}_{g_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2)$ in (A.34), we know $Q_i^*(\boldsymbol{\xi}_{g_i^{(1)}}^{(1)}; \boldsymbol{\xi}_{g^{(2)}}^{(2)}, \mathcal{G}_2) = \sum_j Q_{ij}^*(\Theta_{ij})$ and $Q_i^*(\boldsymbol{\xi}_{g_i^{(1)0}}^{(1)0}; \boldsymbol{\xi}_{g^{(2)0}}^{(2)0}, \mathcal{G}_2^0) = \sum_j Q_{ij}^*(\Theta_{ij}^0)$ by following the similar derivation in (A.1). Since $Q_{ij}(\Theta_{ij}) = \sum_t (\varepsilon_{ij,t} + \mathcal{X}_{ij,t}^\top \Theta_{ij}^0 - \mathcal{X}_{ij,t}^\top \Theta_{ij})^2$, where $\mathcal{X}_{ij,t}$ is defined in (A.31), we can use similar techniques in the proof of Lemma 23 to obtain the conclusion. \blacksquare

Lemma 16. *Under Assumption 4 and A.13, we have $\tau_{\max}^1 = \max_{i,j} \lambda_{\max}(E(\mathcal{X}_{ij,t} \mathcal{X}_{ij,t}^\top)) < \infty$ with $\mathcal{X}_{ij,t}$ defined in (A.31).*

Proof We use the similar techniques in the proof of Lemma 33. Due to the interactive network term, the different term in $\mathcal{X}_{ij,t}$ is $(\mathbf{W}_{\cdot j}^{(2)} \otimes \mathbf{W}_{i \cdot}^{(1)\top})^\top \mathbb{Y}_t$, with $\|\mathbf{W}_{\cdot j}^{(2)} \otimes \mathbf{W}_{i \cdot}^{(1)\top}\|_1 = 1$. Denote $\mathbf{w} = \mathbf{W}_{\cdot j}^{(2)} \otimes \mathbf{W}_{i \cdot}^{(1)\top}$, this typical term has the property that $\text{var}(\mathbf{w}^\top \mathbb{Y}_t) = \mathbf{w}^\top \boldsymbol{\Gamma} \mathbf{w} \leq \|\boldsymbol{\Gamma}\|_{\max} |\mathbf{w}^\top \mathbf{1} \mathbf{1}^\top \mathbf{w}| = \|\boldsymbol{\Gamma}\|_{\max} < c_\Gamma$, where c_Γ is the upper bound conclusion in Lemma 32. Other terms are the same as in the Lemma 33. \blacksquare

Appendix H. Discussion about Future Extension

We remark that the GTNAR model can be easily extended to a general-purpose machine learning model by considering nonlinear terms. Take $q = 2$ for example, model (1) can be

extended to

$$Y_{ij,t} = \lambda_{g_i^{(1)}}^{(1)} \sum_{k=1}^{N_1} \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,(t-1)} + \lambda_{g_j^{(2)}}^{(2)} \sum_{k=1}^{N_2} \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,(t-1)} \quad (\text{A.41})$$

$$+ \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)} + f_1(\mathbf{x}_{it}^{(1)}) + f_2(\mathbf{x}_{jt}^{(2)}) + \varepsilon_{ij,t}, \quad (\text{A.42})$$

where $f_1(\cdot)$ and $f_2(\cdot)$ are non-linear functions. To estimate model (A.42), we could conduct the following (1) and (2) iteratively.

(1) Estimate the linear part. First, given the non-linear part $f(\mathbf{x}_{it}^{(1)})$ and $f(\mathbf{x}_{jt}^{(2)})$, calculate $\tilde{Y}_{ij,t} = Y_{ij,t} - f_1(\mathbf{x}_{it}^{(1)}) - f_2(\mathbf{x}_{jt}^{(2)})$. Then we implement Algorithm A.2 to estimate the other parameters and group memberships iteratively.

(2) Estimate the non-linear part. Given the parameters $\lambda_{g_i^{(1)}}^{(1)}$, $\lambda_{g_j^{(2)}}^{(2)}$, and $\alpha_{g_i^{(1)} g_j^{(2)}}$, calculate $\check{Y}_{ij,t} = Y_{ij,t} - \lambda_{g_i^{(1)}}^{(1)} \sum_{k=1}^{N_1} \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,(t-1)} - \lambda_{g_j^{(2)}}^{(2)} \sum_{k=1}^{N_2} \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,(t-1)} - \alpha_{g_i^{(1)} g_j^{(2)}} Y_{ij,(t-1)}$. Subsequently, one can use deep learning model to fit the function $f_1(\cdot)$ and $f_2(\cdot)$.

Since the discussion of theoretical guarantees is beyond our scope, we leave it as an interesting future topic.

Appendix I. Proof of Main Theorems

I.1 Proof of Theorem 2

Recall that $d(\hat{\Theta}, \Theta)$ is defined in (18). To prove $d(\hat{\Theta}, \Theta^0) = O_p(a_{NT}^2)$, it suffices to show that

$$P\left\{a_{NT}^{-1} \sqrt{d(\hat{\Theta}, \Theta^0)} > K\right\} \rightarrow 0 \text{ as } K \rightarrow \infty.$$

Denote Ω as the parameter space, and partition the space into shells $\Omega_j = \{\Theta \in \Omega : j-1 < a_{NT}^{-1} \sqrt{d(\Theta, \Theta^0)} \leq j\}$, where $j \geq 1$ takes over all positive integers. For a positive integer K , if $a_{NT}^{-1} \sqrt{d(\hat{\Theta}, \Theta^0)} \geq K$, then $\hat{\Theta} \in \Omega_j$ with $j \geq K$. Since $\hat{\Theta}$ minimizes objective function $Q(\Theta)$, we have $(\prod_l N_l T)^{-1} S(\hat{\Theta}) \stackrel{\text{def}}{=} (\prod_l N_l T)^{-1} Q(\hat{\Theta}) - (\prod_l N_l T)^{-1} Q(\Theta^0) \leq 0$. Denote $S^*(\Theta) = E\{S(\Theta)\}$. Then we have

$$\begin{aligned} P\left\{a_{NT}^{-1} \sqrt{d(\Theta, \hat{\Theta})} > K\right\} &\leq P\left\{\inf_{\Theta \in \bigcup_{j \geq K} \Omega_j} \left(\prod_l N_l T\right)^{-1} S(\Theta) \leq 0\right\} \\ &\leq \sum_{j \geq K} P\left\{\inf_{\Theta \in \Omega_j} \left(\prod_l N_l T\right)^{-1} S(\Theta) \leq 0\right\} \\ &\leq \sum_{j \geq K} P\left\{\inf_{\Theta \in \Omega_j} \left(\prod_l N_l T\right)^{-1} \{S(\Theta) - S^*(\Theta)\} \right. \\ &\quad \left. + \inf_{\Theta \in \Omega_j} \left(\prod_l N_l T\right)^{-1} S^*(\Theta) \leq 0\right\} \end{aligned}$$

By Lemma 23, we have that $(\prod_l N_l T)^{-1} S^*(\Theta) = (\prod_l N_l T)^{-1} \{Q^*(\Theta) - Q^*(\Theta^0)\} \geq \tau_{\min} d(\Theta, \Theta^0)$ with τ_{\min} defined in Assumption 2. Since $d(\Theta, \Theta^0) > (j-1)^2 a_{NT}^2$ for all $\Theta \in \Omega_j$, we have

$$\begin{aligned} P\left\{a_{NT}^{-1} \sqrt{d(\Theta, \hat{\Theta})} > K\right\} &\leq \sum_{j \geq K} P\left\{\inf_{\Theta \in \Omega_j} \left(\prod_l N_l T\right)^{-1} \{S(\Theta) - S^*(\Theta)\} \leq -\tau_{\min} (j-1)^2 a_{NT}^2\right\} \\ &\leq \sum_{j \geq K} P\left\{\sup_{\Theta \in \Omega_j} \left(\prod_l N_l T\right)^{-1} |S(\Theta) - S^*(\Theta)| \geq \tau_{\min} (j-1)^2 a_{NT}^2\right\}. \end{aligned} \quad (\text{A.43})$$

Denote $S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) = Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}^0)$, where $Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q})$ is defined in (16) in the main text. For simplicity we denote $d_{i_1 \dots i_q} = d_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}, \Theta_{i_1 \dots i_q}^0)$. We next show that

$$\begin{aligned} &P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \left(\prod_l N_l T\right)^{-1} |S(\Theta) - S^*(\Theta)| \geq x\right\} \\ &\leq C_1 \exp\left\{-C_2 \min(Tx^2/\omega^2, \sqrt{T}x/\omega) + C_3 m + \log\left(\prod_l N_l\right)\right\}. \end{aligned}$$

We have

$$\begin{aligned} &P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \left(\prod_l N_l T\right)^{-1} |S(\Theta) - S^*(\Theta)| \geq x\right\} \\ &\leq P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \left(\prod_l N_l\right)^{-1} \sum_l \sum_{i_l} T^{-1} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})| \geq x\right\} \\ &= P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \frac{1}{\prod_l N_l} \sum_l \sum_{i_l} \frac{T^{-1} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|}{\sqrt{d_{i_1 \dots i_q}}} \sqrt{d_{i_1 \dots i_q}} \geq x\right\} \\ &\leq P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \sqrt{\frac{1}{\prod_l N_l} \sum_l \sum_{i_l} \frac{T^{-2} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|^2}{d_{i_1 \dots i_q}}} \sqrt{d_{i_1 \dots i_q}} \geq x\right\} \\ &= P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \frac{1}{\prod_l N_l} \sum_l \sum_{i_l} \frac{T^{-2} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|^2}{d_{i_1 \dots i_q}} d_{i_1 \dots i_q} \geq x^2\right\} \\ &\leq P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \frac{1}{\prod_l N_l} \sum_l \sum_{i_l} \frac{T^{-2} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|^2}{d_{i_1 \dots i_q}} \omega^2 \geq x^2\right\} \\ &\leq \sum_l \sum_{i_l} P\left\{\sup_{d(\Theta, \Theta^0) \leq \omega^2} \frac{T^{-2} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|^2}{d_{i_1 \dots i_q}} \omega^2 \geq x^2\right\} \\ &\leq \sum_l \sum_{i_l} P\left\{\sup_{\Theta_{i_1 \dots i_q}} \frac{T^{-1} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|}{\sqrt{d_{i_1 \dots i_q}}} \omega \geq x\right\} \\ &\leq C_1 \exp\left\{-C_2 \min(Tx^2/\omega^2, \sqrt{T}x/\omega) + C_3 m + \log\left(\prod_l N_l\right)\right\}, \end{aligned}$$

where the third inequality holds due to $\sum_i (a_i^2 b_i^2) \leq (\sum_i a_i^2)(\sum_i b_i^2)$ and the last inequality holds due to Lemma 19. Here $m = \sum_l (p_l + 1) + 1$ is a fixed constant. Then plugging the last line into (A.43), we have

$$\begin{aligned}
 P\left\{a_{NT}^{-1}\sqrt{d(\boldsymbol{\Theta}, \widehat{\boldsymbol{\Theta}})} > K\right\} &\leq \sum_{j \geq K} P\left\{\sup_{\boldsymbol{\Theta} \in \Omega_j} \left(\prod_l N_l T\right)^{-1} |S(\boldsymbol{\Theta}) - S^*(\boldsymbol{\Theta})| \geq \tau_{\min}(j-1)^2 a_{NT}^2\right\} \\
 &\leq \sum_{j \geq K} C_1 \exp\left\{-C_2 \min\left(\frac{T(j-1)^4 a_{NT}^2}{j^2}, \frac{\sqrt{T}(j-1)^2 a_{NT}}{j}\right) + C_3 m + \log\left(\prod_l N_l\right)\right\} \\
 &\leq \sum_{j \geq K} C_1 \exp\left\{-C_2 \min\left(\frac{(j-1)^4 \log^2(\prod_l N_l)}{j^2}, \frac{(j-1)^2 \log(\prod_l N_l)}{j}\right) + C_3 m + \log\left(\prod_l N_l\right)\right\} \\
 &= \sum_{j \geq K} C_1 \exp\left\{-C_2 \frac{(j-1)^2 \log(\prod_l N_l)}{j} + C_3 m + \log\left(\prod_l N_l\right)\right\} \\
 &\leq \sum_{j \geq K} C_1 \exp\left\{-C_2(j-1)(1-K^{-1})\log\left(\prod_l N_l\right) + C_3 m + \log\left(\prod_l N_l\right)\right\} \\
 &= \frac{\exp\left\{-C_2(K-1)(1-K^{-1})\log(\prod_l N_l) + C_3 m + \log(\prod_l N_l)\right\}}{1 - \exp\left\{-C_2(1-K^{-1})\log(\prod_l N_l)\right\}} \\
 &\leq \frac{\exp(-C_2(K-1)(1-K^{-1}) - 1)\log(\prod_l N_l) + C_3 m}{1 - \exp(-C_2 \log(\prod_l N_l))} \rightarrow 0
 \end{aligned} \tag{A.44}$$

when $K \rightarrow \infty$. The third inequality holds by plugging in the expression of $a_{NT}^2 = \log^2(\prod_l N_l) T^{-1}$. The equality holds due to the sum of a geometric sequence. Then we finish the proof.

1.2 Proof of Theorem 3

In the following proof we write $Q(\widehat{\boldsymbol{\xi}}(\underline{G}), \widehat{\boldsymbol{G}}(\underline{G}))$ as $Q(\widehat{\boldsymbol{\Theta}}(\underline{G}))$ as we defined in (16). We use the notations here to denote $\widehat{\boldsymbol{\Theta}}(\underline{G})$ as the estimate when group numbers G_1, \dots, G_q are specified. Define $G_{-l} = (G_k : k \neq l)^\top \in \mathbb{R}^{q-1}$ and we use $G_{-l} \geq G_{-l,0}$ to denote that $G_k \geq G_{k,0}$ for $k \neq l$. To establish the consistency property, we consider both the underfitted model (there exists an l with $G_l < G_{l,0}$) and the overfitted model ($\{G_l > G_{l,0}, G_{-l} \geq G_{-l,0}\}$ for any l). In the following we show $\text{QIC}(\underline{G}_0) < \text{QIC}(\underline{G})$ with probability tending to 1 under both circumstances. We write the difference of the two criteria as

$$\begin{aligned}
 \text{QIC}(\underline{G}) - \text{QIC}(\underline{G}_0) &= \log\{Q(\widehat{\boldsymbol{\Theta}}(\underline{G}))\} - \log\{Q(\widehat{\boldsymbol{\Theta}}(\underline{G}_0))\} + \lambda(\underline{G}) - \lambda(\underline{G}_0) \\
 &= \log\left\{1 + \frac{Q(\widehat{\boldsymbol{\Theta}}(\underline{G})) - Q(\widehat{\boldsymbol{\Theta}}(\underline{G}_0))}{Q(\widehat{\boldsymbol{\Theta}}(\underline{G}_0))}\right\} + \lambda(\underline{G}) - \lambda(\underline{G}_0).
 \end{aligned}$$

1. **OVERFITTED MODEL.** In the following we show that $(\prod_l N_l T)^{-1}\{Q(\widehat{\boldsymbol{\Theta}}(\underline{G}_0)) - Q(\widehat{\boldsymbol{\Theta}}(\underline{G}))\} = (\prod_l N_l T)^{-1}\{[Q(\widehat{\boldsymbol{\Theta}}(\underline{G}_0)) - Q(\boldsymbol{\Theta}^0)] - [Q(\widehat{\boldsymbol{\Theta}}(\underline{G})) - Q(\boldsymbol{\Theta}^0)]\} = (\prod_l N_l T)^{-1}\{S(\widehat{\boldsymbol{\Theta}}(\underline{G}_0)) - S(\widehat{\boldsymbol{\Theta}}(\underline{G}))\} = O_p(a_{NT})$, where $a_{NT} = T^{-1}(\sum_l \log N_l)^2$. Note that we have

$$\frac{1}{(\prod_l N_l)T} |S(\widehat{\boldsymbol{\Theta}}(\underline{G})) - S(\widehat{\boldsymbol{\Theta}}(\underline{G}_0))| \leq \frac{1}{(\prod_l N_l)T} |S(\widehat{\boldsymbol{\Theta}}(\underline{G}_0)) - S^*(\widehat{\boldsymbol{\Theta}}(\underline{G}_0))|$$

$$\begin{aligned}
 & + \frac{1}{(\prod_l N_l)T} |S^*(\widehat{\Theta}(\underline{G}_0)) - S^*(\widehat{\Theta}(\underline{G}))| + \frac{1}{(\prod_l N_l)T} |S(\widehat{\Theta}(\underline{G})) - S^*(\widehat{\Theta}(\underline{G}))| \\
 & = O_p(T^{-1}(\sum_l \log N_l)^2) + \frac{1}{(\prod_l N_l)T} |S^*(\widehat{\Theta}(\underline{G}_0)) - S^*(\widehat{\Theta}(\underline{G}))|
 \end{aligned}$$

by Lemma 19.

Further by Lemma 23 and Theorem 2, it holds that

$$\begin{aligned}
 & \frac{1}{(\prod_l N_l)T} |S^*(\widehat{\Theta}(\underline{G}_0)) - S^*(\widehat{\Theta}(\underline{G}))| \\
 & \leq \frac{1}{(\prod_l N_l)T} |Q^*(\widehat{\Theta}(\underline{G}_0)) - Q^*(\Theta^0)| + \frac{1}{(\prod_l N_l)T} |Q^*(\Theta^0) - Q^*(\widehat{\Theta}(\underline{G}))| \\
 & \leq \tau_{\max} \left\{ d(\widehat{\Theta}(\underline{G}_0), \Theta^0) + d(\widehat{\Theta}(\underline{G}), \Theta^0) \right\} = O_p(T^{-1}(\sum_l \log N_l)^2),
 \end{aligned}$$

where τ_{\max} is defined in Assumption 2. Since m is a fixed constant, $((\prod_l N_l)T)^{-1} |Q(\widehat{\Theta}(\underline{G}_0)) - Q(\widehat{\Theta}(\underline{G}))| = O_p\{T^{-1}(\sum_l \log N_l)^2\}$, i.e.,

$$\frac{1}{(\prod_l N_l)T} Q(\widehat{\Theta}(\underline{G})) = \frac{1}{(\prod_l N_l)T} Q(\widehat{\Theta}(\underline{G}_0)) + O_p\{T^{-1}(\sum_l \log N_l)^2\}.$$

Next we show $((\prod_l N_l)T)^{-1} Q(\widehat{\Theta}(\underline{G}_0)) = \sigma^2 + o_p(1)$. Similar to above while replacing $\widehat{\Theta}(\underline{G})$ by Θ^0 , we have

$$((\prod_l N_l)T)^{-1} Q(\widehat{\Theta}(\underline{G}_0)) = ((\prod_l N_l)T)^{-1} Q^*(\Theta^0) + O_p\{T^{-1}(\sum_l \log N_l)^2\}. \quad (\text{A.45})$$

Here we have $((\prod_l N_l)T)^{-1} Q^*(\Theta^0) = \sigma^2 > 0$. This shows that $((\prod_l N_l)T)^{-1} Q(\widehat{\Theta}(\underline{G}_0)) = \sigma^2 + o_p(1)$. It implies $\text{QIC}(\underline{G}) - \text{QIC}(\underline{G}_0) = O_p\{T^{-1}(\sum_l \log N_l)^2\} + \lambda(\underline{G}) - \lambda(\underline{G}_0)$. We can verify that $\lambda(\underline{G}) - \lambda(\underline{G}_0) \geq \kappa$ under the overfitting case, where $\kappa = \lambda(\underline{G})/(\sum_l G_l)$. Further note that $\kappa \gg T^{-1}(\sum_l \log N_l)^2$ as we assume, this implies $\text{QIC}(\underline{G}_0) < \text{QIC}(\underline{G})$ with probability tending to 1 under the overfitting case.

2. UNDERFITTED MODEL. From the specification of $\lambda(\underline{G})$, we have $|\lambda(\underline{G}) - \lambda(\underline{G}_0)| = o(c_{\text{gap}}/(\prod_l G_l))$. It suffices to show that $((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G})) - Q(\widehat{\Theta}(\underline{G}_0))\} \geq C_1 c_{\text{gap}}(c_\pi)^q / (\prod_l G_l)$ for a positive constant C_1 when $N_l \rightarrow \infty$. Note that

$$\begin{aligned}
 & ((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G})) - Q(\widehat{\Theta}(\underline{G}_0))\} \\
 & = ((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G})) - Q^*(\Theta^0) - Q(\widehat{\Theta}(\underline{G}_0)) + Q^*(\Theta^0)\},
 \end{aligned}$$

and we have already proved that $((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G}_0)) - Q^*(\Theta^0)\} = O_p(T^{-1}(\sum_l \log N_l)^2) = o_p\{c_{\text{gap}}/(\prod_l G_l)\}$ in (A.45). Also note that c_π is a constant according to Assumption 7. Hence it suffices to show $((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G})) - Q^*(\Theta^0)\} \geq C c_{\text{gap}}(c_\pi)^q / (\prod_l G_l)$ for some positive constant C when $N_l \rightarrow \infty$.

Without loss of generality, we consider $G_l < G_{l,0}$. To prove the result, we use two steps. First, we show that by Assumption 6, we have

$$\begin{aligned} \max_{g^{(l)} \in [G_{l,0}]} \max_{g^{-(l)} \in [G_{-l}^0]} \min_{g^{(l)'} \in [G_l]} \left\{ \|\widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'}} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} \right. \\ \left. - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0 \right\} \geq c_{\text{gap}}/4, \quad (\text{A.46}) \end{aligned}$$

where $\varphi_l(g^{(l)}) = \operatorname{argmax}_{g^{(l)'} \in [G_l]} \sum_{i_l} I(g_{i_l}^{(l)0} = g^{(l)}, \widehat{g}_{i_l}^{(l)} = g^{(l)'})$, and $\{g^{-(l)} \in [G_{-l}^0]\}$ denotes $\{g^{(1)} \in [G_{1,0}], \dots, g^{(l-1)} \in [G_{l-1,0}], g^{(l+1)} \in [G_{l+1,0}], \dots, g^{(q)} \in [G_{q,0}]\}$. We prove (A.46) by contradiction. Assume (A.46) does not hold. Define $\sigma_l : [G_{l,0}] \rightarrow [G_l]$, such that

$$\begin{aligned} \sigma_l(g^{(l)}) = \operatorname{argmin}_{g^{(l)' \in [G_l]} \left\{ \|\widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 \right. \\ \left. + \max_{g^{-(l)} \in [G_{-l}^0]} |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'}} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} \right. \\ \left. - \alpha_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})}^0 \right\}. \end{aligned}$$

Since $G_{l,0} > G_l$, there exists at least two $g_1^{(l)}, g_2^{(l)} \in [G_{l,0}]$, such that $\sigma_l(g_1^{(l)}) = \sigma_l(g_2^{(l)})$. Then we have for all $g^{-(l)} \in [G_{-l}^0]$,

$$\begin{aligned} & \|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_2^{(l)}}^{(l)0}\|^2 + |\alpha_{g^{(1)} \dots g^{(l-1)} g_1^{(l)} g^{(l+1)} \dots g^{(q)}}^0 - \alpha_{g^{(1)} \dots g^{(l-1)} g_2^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \\ & \leq 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{\sigma_l(g_1^{(l)})}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{\sigma_l(g_2^{(l)})}^{(l)}\|^2 \\ & + 2\left\{ |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) \sigma_l(g_1^{(l)}) \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g_1^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right. \\ & \left. + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) \sigma_l(g_2^{(l)}) \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g_2^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\} \\ & < c_{\text{gap}}/2 + c_{\text{gap}}/2 = c_{\text{gap}}, \quad (\text{A.47}) \end{aligned}$$

which contradicts Assumption 6. We explain the last inequality as follows. For any $g^{(l)} \in [G_{l,0}]$, $g^{-(l)} \in [G_{-l}^0]$, $g^{(l)' \in [G_l]$, we have

$$\begin{aligned} & \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) \sigma_l(g^{(l)}) \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \\ & \leq \max_{g^{-(l)} \in [G_{-l}^0]} \left\{ \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)}\|^2 \right. \\ & \left. + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) \sigma_l(g^{(l)}) \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\} \\ & \leq \max_{g^{-(l)} \in [G_{-l}^0]} \left\{ \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)}\|^2 \right. \\ & \left. + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'}} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\} \\ & \leq \max_{g^{(l)} \in [G_{l,0}]} \max_{g^{-(l)} \in [G_{-l}^0]} \left\{ \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)}\|^2 \right. \end{aligned}$$

$$+ |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \},$$

where the second inequality holds by the definition of $\sigma_l(g^{(l)})$. Since the above inequality holds for any $g^{(l)'}$, hence by setting

$$g^{(l)'} = \operatorname{argmin}_{g^{(l)'} \in [G_l]} \left\{ \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\},$$

we get

$$\begin{aligned} & \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) \sigma_l(g^{(l)}) \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \\ & \leq \max_{g^{(l)} \in [G_{l,0}]} \max_{g^{-(l)} \in [G_{-l}^0]} \min_{g^{(l)'} \in [G_l]} \left\{ \|\boldsymbol{\theta}_{g^{(l)}}^{(l)0} - \widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\} < c_{\text{gap}}/4. \end{aligned}$$

Hence, by the contradiction (A.47), we could obtain that (A.46) holds.

Subsequently, define the mapping $\varphi_l(g^{(l)}) = \operatorname{argmax}_{g^{(l)'} \in [G_l]} \sum_{i_l} I(g_{i_l}^{(l)0} = g^{(l)}, \widehat{g}_{i_l}^{(l)} = g^{(l)'})$. In addition, let

$$\begin{aligned} (g^{(1)*}, \dots, g^{(q)*}) &= \operatorname{argmax}_{g^{(1)} \in [G_{1,0}], \dots, g^{(q)} \in [G_{q,0}]} \left[\min_{g^{(l)'} \in [G_l]} \left\{ \|\widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 \right. \right. \\ & \left. \left. + |\widehat{\alpha}_{\varphi_1(g^{(1)}) \dots \varphi_{l-1}(g^{(l-1)}) g^{(l)'} \varphi_{l+1}(g^{(l+1)}) \dots \varphi_q(g^{(q)})} - \alpha_{g^{(1)} \dots g^{(l-1)} g^{(l)} g^{(l+1)} \dots g^{(q)}}^0|^2 \right\} \right]. \quad (\text{A.48}) \end{aligned}$$

In this case, we have

$$\begin{aligned} d(\widehat{\boldsymbol{\Theta}}(\underline{G}), \boldsymbol{\Theta}^0) &= \frac{1}{\prod_l N_l} \sum_{i_1=1}^{N_1} \dots \sum_{i_q=1}^{N_q} \left\| \widehat{\boldsymbol{\Theta}}_{i_1 \dots i_q} - \boldsymbol{\Theta}_{i_1 \dots i_q} \right\|^2 \\ &= \sum_l \frac{1}{N_l} \sum_{i_l=1}^{N_l} \sum_{g^{(l)'}=1}^{G_l} \sum_{g^{(l)}=1}^{G_{l,0}} I(\widehat{g}_{i_l}^{(l)} = g^{(l)'}, g_{i_l}^{(l)0} = g^{(l)}) \|\widehat{\boldsymbol{\theta}}_{g^{(l)'}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 \\ &+ \frac{1}{\prod_l N_l} \sum_{i_1, \dots, i_q} \sum_{g^{(1)'}=1}^{G_1} \dots \sum_{g^{(q)'}=1}^{G_q} \sum_{g^{(1)}=1}^{G_{1,0}} \dots \sum_{g^{(q)}=1}^{G_{q,0}} \\ & I(\widehat{g}_{i_1}^{(1)} = g^{(1)'}, \dots, \widehat{g}_{i_q}^{(q)} = g^{(q)'}, g_{i_1}^{(1)0} = g^{(1)}, \dots, g_{i_q}^{(q)0} = g^{(q)}) |\widehat{\alpha}_{g^{(1)'} \dots g^{(q)'}} - \alpha_{g^{(1)} \dots g^{(q)}}^0|^2 \\ &\geq \sum_l \frac{1}{N_l} \sum_{i_l=1}^{N_l} I(\widehat{g}_{i_l}^{(l)} = \varphi_l(g^{(l)*}), g_{i_l}^{(l)0} = g^{(l)*}) \|\widehat{\boldsymbol{\theta}}_{\varphi_l(g^{(l)*})}^{(l)} - \boldsymbol{\theta}_{g^{(l)*}}^{(l)0}\|^2 \\ &+ \frac{1}{\prod_l N_l} \sum_{i_1, \dots, i_q} I(\widehat{g}_{i_1}^{(1)} = \varphi_1(g^{(1)*}), \dots, \widehat{g}_{i_q}^{(q)} = \varphi_q(g^{(q)*}), g_{i_1}^{(1)0} = g^{(1)*}, \dots, g_{i_q}^{(q)0} = g^{(q)*}) \\ & |\widehat{\alpha}_{\varphi_1(g^{(1)*}) \dots \varphi_q(g^{(q)*})} - \alpha_{g^{(1)*} \dots g^{(q)*}}^0|^2. \end{aligned}$$

By (A.46), we further have

$$\begin{aligned}
 & d(\widehat{\Theta}(\underline{G}), \Theta^0) \\
 & \geq \frac{1}{\prod_l N_l} \sum_{i_1, \dots, i_q} I(\widehat{g}_{i_1}^{(1)} = \varphi_1(g^{(1)*}), \dots, \widehat{g}_{i_q}^{(q)} = \varphi_q(g^{(q)*}), g_{i_1}^{(1)0} = g^{(1)*}, \dots, g_{i_q}^{(q)0} = g^{(q)*}) \\
 & \quad \|\widehat{\theta}_{\varphi_l(g^{(l)*})}^{(l)} - \theta_{g^{(l)*}}^{(l)0}\|^2 \\
 & + \frac{1}{\prod_l N_l} \sum_{i_1, \dots, i_q} I(\widehat{g}_{i_1}^{(1)} = \varphi_1(g^{(1)*}), \dots, \widehat{g}_{i_q}^{(q)} = \varphi_q(g^{(q)*}), g_{i_1}^{(1)0} = g^{(1)*}, \dots, g_{i_q}^{(q)0} = g^{(q)*}) \\
 & \quad |\widehat{\alpha}_{\varphi_1(g^{(1)*}) \dots \varphi_q(g^{(q)*})} - \alpha_{g^{(1)*} \dots g^{(q)*}}^0|^2 \\
 & \geq \frac{c_{\text{gap}}}{4 \prod_l N_l} \sum_{i_1, \dots, i_q} I(\widehat{g}_{i_1}^{(1)} = \varphi_1(g^{(1)*}), \dots, \widehat{g}_{i_q}^{(q)} = \varphi_q(g^{(q)*}), g_{i_1}^{(1)0} = g^{(1)*}, \dots, g_{i_q}^{(q)0} = g^{(q)*}) \\
 & = \frac{c_{\text{gap}}}{4 \prod_l N_l} \sum_{l=1}^q \sum_{i_l=1}^{N_l} I\{\widehat{g}_{i_l}^{(l)} = \varphi_l(g^{(l)*}), g_{i_l}^{(l)0} = g^{(l)*}\}.
 \end{aligned}$$

where the last inequality is due to that

$$\begin{aligned}
 & \|\widehat{\theta}_{\varphi_l(g^{(l)*})}^{(l)} - \theta_{g^{(l)*}}^{(l)0}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)*}) \dots \varphi_q(g^{(q)*})} - \alpha_{g^{(1)*} \dots g^{(q)*}}^0|^2 \\
 & \geq \min_{g^{(l)'} \in [G_l]} \left\{ \|\widehat{\theta}_{g^{(l)'}}^{(l)} - \theta_{g^{(l)*}}^{(l)0}\|^2 + |\widehat{\alpha}_{\varphi_1(g^{(1)*}) \dots \varphi_{l-1}(g^{(l-1)*}) g^{(l)'} \varphi_{l+1}(g^{(l+1)*}) \dots \varphi_q(g^{(q)*})} - \alpha_{g^{(1)*} \dots g^{(q)*}}^0|^2 \right\} \\
 & = \max_{g^{(l)} \in [G_{l,0}]} \max_{g^{-(l)} \in [G_{-l}^0]} \min_{g^{(l)'} \in [G_l]} \left\{ \|\widehat{\theta}_{g^{(l)'}}^{(l)} - \theta_{g^{(l)}}^{(l)0}\|^2 \right. \\
 & \quad \left. + |\widehat{\alpha}_{\varphi_1(g^{(1)*}) \dots \varphi_{l-1}(g^{(l-1)*}) g^{(l)'} \varphi_{l+1}(g^{(l+1)*}) \dots \varphi_q(g^{(q)*})} - \alpha_{g^{(1)*} \dots g^{(q)*}}^0|^2 \right\} \geq \frac{c_{\text{gap}}}{4}
 \end{aligned}$$

by the definition (A.48) and (A.46).

By the definition of $\varphi_l(\cdot)$, we have

$$N_l^{-1} \sum_{i_l} I(\widehat{g}_{i_l}^{(l)} = \varphi_l(g^{(l)*}), g_{i_l}^{(l)0} = g^{(l)*}) \geq N_l^{-1} G_l^{-1} \sum_{i_l} I(g_{i_l}^{(l)0} = g^{(l)*}) = G_l^{-1} \pi_{g^{(l)*}, N_l}^{(l)}$$

From Assumption 7, we know that there exists $M > 0$, such that for $N_l > M$ it holds $G_l^{-1} \pi_{g^{(l)*}, N_l}^{(l)} > c_\pi / 2G_l$. This yields $d(\widehat{\Theta}(\underline{G}), \Theta^0) \geq \frac{c_{\text{gap}}(c_\pi)^q}{2^{q+2} \prod_l G_l}$ as $N_l \rightarrow \infty$ for all $l \in [q]$ by Assumption 7. Then by Lemmas 19 and 23, we have

$$\begin{aligned}
 & ((\prod_l N_l)T)^{-1} \{Q(\widehat{\Theta}(\underline{G})) - Q(\widehat{\Theta}(\underline{G}_0))\} \\
 & = O_p(T^{-1}(\sum_l \log N_l)^2) + ((\prod_l N_l)T)^{-1} \{Q^*(\widehat{\Theta}(\underline{G})) - Q^*(\Theta^0)\} \\
 & \geq \tau_{\min} d(\widehat{\Theta}(\underline{G}), \Theta^0) + O_p(T^{-1}(\sum_l \log N_l)^2) \\
 & \geq \frac{\tau_{\min} c_{\text{gap}} (c_\pi)^q}{2^{q+2} \prod_l G_l} + O_p(T^{-1/2}(m + \sum_l \log N_l)) \geq \frac{C c_{\text{gap}} (c_\pi)^q}{\prod_l G_l}.
 \end{aligned}$$

with probability tending to 1, where $C = \tau_{\min} / 2^{q+2}$ is a positive constant, and the last inequality holds since the condition (19) that $T^{-1}(\sum_l \log N_l)^2 \ll c_{\text{gap}} / (\prod_l G_l)$.

I.3 Proof of Proposition 4

Recall that

$$Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-l}}^{-(l)}, \mathcal{G}_{-l}) = \sum_{m \neq l} \sum_{i_m=1}^{N_m} \sum_{t=1}^T \left\{ Y_{i_1 \dots i_q, t} - \sum_{l=1}^q \lambda_{g_{i_l}^{(l)}}^{(l)} \sum_{k=1}^{N_l} w_{i_l k}^{(l)} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)} \right. \\ \left. - \alpha_{g_{i_1}^{(1)} \dots g_{i_q}^{(q)}} Y_{i_1 \dots i_q, (t-1)} - \sum_{l=1}^q \mathbf{x}_{i_l t}^{(l)\top} \boldsymbol{\zeta}_{g_{i_l}^{(l)}}^{(l)} \right\}^2,$$

which is defined in (A.1). Define

$$\sigma_l(g^{(l)}) = \operatorname{argmin}_{\widehat{g}^{(l)} \in [G_l]} \left\{ \|\widehat{\boldsymbol{\theta}}_{\widehat{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 \right. \\ \left. + \frac{1}{\prod_{m \neq l} N_m} \sum_{i_{-l}} |\widehat{\alpha}_{g_{i_1}^{(1)} \dots g_{i_{l-1}}^{(l-1)} \widehat{g}^{(l)} g_{i_{l+1}}^{(l+1)} \dots g_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)0} \dots g_{i_{l-1}}^{(l-1)0} g^{(l)} g_{i_{l+1}}^{(l+1)0} \dots g_{i_q}^{(q)0}}|^2 \right\}.$$

Let $R_1 = \sup_l \sup_{i_l} \sup_{\|\boldsymbol{\Theta}_{i_1 \dots i_q}\|_{\max} < R} \frac{1}{T} \{ |Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) - Q_{i_1 \dots i_q}^*(\boldsymbol{\Theta}_{i_1 \dots i_q})| \}$ and we have $R_1 = O_p(T^{-1/2}(m + \sum_l \log N_l))$ due to Lemma 17. Therefore again by Lemma 17 and the fact that $\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} = (\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)\top}, \operatorname{vec}(\widehat{\boldsymbol{\alpha}}_{\widehat{g}_{i_l}^{(l)}}))^\top$ minimize $Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-l}}^{-(l)}, \mathcal{G}_{-l})$, we have

$$\frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - R_1 \leq \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) \\ \leq \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) \leq \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) + R_1$$

for all $i_l \in [N_l]$. Therefore it holds

$$\sup_{i_l} \left\{ \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) \right\} \leq 2R_1. \quad (\text{A.49})$$

On the other hand by Lemma 23 we have

$$0 \leq \sup_{i_l} \left\{ \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-l}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)0}; \boldsymbol{\xi}_{g^{-l}}^{-(l)0}, \mathcal{G}_{-l}^0) \right\} \\ \leq \tau_{\max} \sup_{i_l} \left\{ \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_m}^{(m)}}^{(m)} - \boldsymbol{\theta}_{g_{i_m}^{(m)0}}^{(m)0}\|^2 + \|\widehat{\boldsymbol{\theta}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 \right. \\ \left. + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m=1}^{N_m} |\widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_{l-1}}^{(l-1)} \sigma_l(g_{i_l}^{(l)0}) \widehat{g}_{i_{l+1}}^{(l+1)} \dots \widehat{g}_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)0} \dots g_{i_q}^{(q)0}}|^2 \right\} \\ = O_p(T^{-1}(\sum_l \log N_l)^2) + O_p(T^{-1}(\sum_l \log N_l)^2) = O_p(T^{-1}(\sum_l \log N_l)^2) \quad (\text{A.50})$$

where the second last equation is obtained by Theorem 2 and Lemma 24.

Next, we use the order in (A.50) to improve the order of R_1 in (A.49). Denote $\widehat{\Theta}_{i_l}^\sigma$ as the parameter by replacing the $\widehat{g}_{i_l}^{(l)}$ by $\sigma(g_{i_l}^{(l)0})$ in $\widehat{\Theta}_{i_l}$. while keep the other terms the same. By (A.50), there exists a positive constant C , the event $\mathcal{O}_1 = \{\sup_{i_l} d_{i_l}(\widehat{\Theta}_{i_l}^\sigma, \Theta_{i_l}^0) \leq C a_{NT}^2\}$ has the probability $P(\mathcal{O}_1) > 1 - \epsilon$ for $\forall \epsilon > 0$, where $a_{NT} = (\sum_l \log N_l) / \sqrt{T}$. To derive the order of R_1 , we need to derive the order of distance $d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma)$. We next show that $P(a_{NT}^{-1} \sup_{i_l} \sqrt{d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma)} > K) \rightarrow 0$ as $K \rightarrow \infty$, which implies that $d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma) = O_p(a_{NT}^2)$. Partition the parameter Θ_{i_l} into space $\Omega_j = \{\Theta_{i_l} \in \Omega : \gamma(j-1) < a_{NT}^{-1} \sqrt{d_{i_l}(\Theta_{i_l}, \widehat{\Theta}_{i_l}^\sigma)} \leq \gamma j\}$, with $j \geq 1$ taking over all positive integers and $\gamma = \sqrt{1 + C(\tau_{\max} + \tau_{\min})} / \tau_{\min}$ is a positive constant. For a positive integer K , if $a_{NT}^{-1} \sqrt{d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma)} > K$, then $\widehat{\Theta}_{i_l} \in \bigcup_{j \geq K} \Omega_j$. Define

$$S_{i_l}(\Theta_{i_l}) = Q_{i_l}(\boldsymbol{\xi}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - Q_{i_l}(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l}).$$

By the definition of $\widehat{\Theta}_{i_l}$, we know that $S_{i_l}(\widehat{\Theta}_{i_l}) \leq 0$. Hence, we have

$$\begin{aligned} & P\left\{a_{NT}^{-1} \sup_{i_l} \sqrt{d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma)} \geq K\right\} \leq P\left\{\inf_{i_l} \inf_{\Theta_{i_l} \in \bigcup_{j \geq K} \Omega_j} \left(\prod_{m \neq l} N_m T\right)^{-1} S_{i_l}(\widehat{\Theta}_{i_l}) \leq 0\right\} \\ & \leq \sum_{j \geq K} P\left\{\inf_{i_l} \left[\inf_{\Theta_{i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T\right)^{-1} \{S_{i_l}(\widehat{\Theta}_{i_l}) - S_{i_l}^*(\widehat{\Theta}_{i_l})\}\right.\right. \\ & \quad \left.\left.+ \inf_{\Theta_{i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T\right)^{-1} S_{i_l}^*(\widehat{\Theta}_{i_l})\right] \leq 0\right\} \end{aligned}$$

Under event \mathcal{O}_1 , we know that if $\widehat{\Theta}_{i_l} \in \Omega_j$,

$$\begin{aligned} & \left(\prod_{m \neq l} N_m T\right)^{-1} \inf_{i_l} S_{i_l}^*(\widehat{\Theta}_{i_l}) = \left(\prod_{m \neq l} N_m T\right)^{-1} \inf_{i_l} \{Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l})\} \\ & = \left(\prod_{m \neq l} N_m T\right)^{-1} \inf_{i_l} \left[\{Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}^{(l)}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l}) - Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)}; \boldsymbol{\xi}_{g_{i_l}^{(l)}}^{-(l)}, \mathcal{G}_{-l}^0)\}\right. \\ & \quad \left.+ \{Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)}; \boldsymbol{\xi}_{g_{i_l}^{(l)}}^{-(l)}, \mathcal{G}_{-l}^0) - Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \widehat{\boldsymbol{\xi}}_{g_{i_l}^{(l)}}^{-(l)}, \widehat{\mathcal{G}}_{-l})\}\right] \\ & \stackrel{\textcircled{1}}{\geq} \tau_{\min} \inf_{i_l} d_{i_l}(\widehat{\Theta}_{i_l}, \Theta_{i_l}^0) - \tau_{\max} \sup_{i_l} d_{i_l}(\widehat{\Theta}_{i_l}^\sigma, \Theta_{i_l}^0) \\ & \stackrel{\textcircled{2}}{\geq} \tau_{\min} \inf_{i_l} d_{i_l}(\widehat{\Theta}_{i_l}, \widehat{\Theta}_{i_l}^\sigma) - (\tau_{\min} + \tau_{\max}) \sup_{i_l} d_{i_l}(\widehat{\Theta}_{i_l}^\sigma, \Theta_{i_l}^0) \\ & \stackrel{\textcircled{3}}{\geq} \tau_{\min} (j-1)^2 (1 + C \frac{\tau_{\max} + \tau_{\min}}{\tau_{\min}}) a_{NT}^2 - C(\tau_{\max} + \tau_{\min}) a_{NT}^2 \\ & = \tau_{\min} a_{NT}^2 (j-1)^2 + C\{(j-1)^2 - 1\} (\tau_{\min} + \tau_{\max}) a_{NT}^2 \\ & \stackrel{\textcircled{4}}{\geq} \tau_{\min} a_{NT}^2 (j-1)^2, \end{aligned}$$

where the inequality ① holds due to Lemma 23, inequality ② holds due to the triangle inequality, inequality ③ holds under \mathcal{O}_1 and the definition of the partition, and inequality ④ holds when $j > 1$. Therefore, we have

$$\begin{aligned}
 & P\left\{a_{NT}^{-1} \sup_{i_l} \sqrt{d_{i_l}(\widehat{\Theta}_{\cdot i_l}, \widehat{\Theta}_{\cdot i_l}^\sigma)} \geq K\right\} \\
 & \leq \sum_{j \geq K} P\left\{\inf_{i_l} \left[\inf_{\widehat{\Theta}_{\cdot i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T \right)^{-1} \{S_{i_l}(\widehat{\Theta}_{\cdot i_l}) - S_{i_l}^*(\widehat{\Theta}_{\cdot i_l})\} \right. \right. \\
 & \quad \left. \left. + \inf_{\widehat{\Theta}_{\cdot i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T \right)^{-1} S_{i_l}^*(\widehat{\Theta}_{\cdot i_l}) \right] \leq 0\right\} \\
 & \leq \sum_{j \geq K} P\left\{\sup_{i_l} \sup_{\widehat{\Theta}_{\cdot i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T \right)^{-1} |S_{i_l}(\widehat{\Theta}_{\cdot i_l}) - S_{i_l}^*(\widehat{\Theta}_{\cdot i_l})| \geq \tau_{\min} a_{NT}^2 (j-1)^2\right\} + P(\mathcal{O}_1^c).
 \end{aligned} \tag{A.51}$$

Next, we derive the upper bound for the first term in (A.51) under event \mathcal{O}_1 . Note that

$$\begin{aligned}
 S_{i_l}(\widehat{\Theta}_{\cdot i_l}) &= Q_{i_l}(\widehat{\Theta}_{\cdot i_l}) - Q_{i_l}(\widehat{\Theta}_{\cdot i_l}^\sigma) \\
 &= \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ \left(Y_{i_1 \dots i_q, t} - \mathcal{X}_{i_1 \dots i_q, t}^\top \widehat{\Theta}_{i_1 \dots i_q} \right)^2 - \left(Y_{i_1 \dots i_q, t} - \mathcal{X}_{i_1 \dots i_q, t}^\top \widehat{\Theta}_{i_1 \dots i_q}^\sigma \right)^2 \right\} \\
 &= \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q})^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}) \right. \\
 & \quad \left. - (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma) \right. \\
 & \quad \left. + 2 \mathcal{X}_{i_1 \dots i_q, t}^\top (\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q}) \varepsilon_{i_1 \dots i_q, t} \right\}
 \end{aligned}$$

Hence we have

$$\begin{aligned}
 & |S_{i_l}(\widehat{\Theta}_{\cdot i_l}) - S_{i_l}^*(\widehat{\Theta}_{\cdot i_l})| \leq \\
 & \leq 2 \left| \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma) \right. \right. \\
 & \quad \left. \left. - [(\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)] \right\} \right| \tag{A.52}
 \end{aligned}$$

$$\begin{aligned}
 & + \left| \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ (\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q})^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q}) \right. \right. \\
 & \quad \left. \left. - [(\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q})^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) (\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q})] \right\} \right| \tag{A.53}
 \end{aligned}$$

$$\begin{aligned}
 & + 2 \left| \sum_{m \neq l} \sum_{i_m} \sum_t \mathcal{X}_{i_1 \dots i_q, t}^\top (\widehat{\Theta}_{i_1 \dots i_q}^\sigma - \widehat{\Theta}_{i_1 \dots i_q}) \varepsilon_{i_1 \dots i_q, t} \right| \tag{A.54}
 \end{aligned}$$

Due to the definition of parameter Θ_{i_l} 's partition Ω_j and under event \mathcal{O}_1 , the concentration inequalities for (A.52)–(A.54) can be obtained by Lemma 20 that,

$$\sum_{j \geq K} P\left\{\sup_{i_l} \sup_{\widehat{\Theta}_{\cdot i_l} \in \Omega_j} \left(\prod_{m \neq l} N_m T \right)^{-1} \right.$$

$$\begin{aligned}
 & \left| \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma) \right. \right. \\
 & \quad \left. \left. - E[(\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \widehat{\Theta}_{i_1 \dots i_q}^\sigma)] \right\} \right| \geq \tau_{\min} a_{NT}^2 (j-1)^2 \Big\} \\
 & \leq \sum_{j \geq K} C_1 \exp \left\{ -C_2 \min \left(\frac{T(j-1)^4 a_{NT}^2}{C}, \frac{\sqrt{T}(j-1)^2 a_{NT}}{\sqrt{C}} \right) + C_3 m + \log \left(\prod_l N_l \right) \right\} \\
 & = \sum_{j \geq K} C_1 \exp \left\{ -\widetilde{C}_2 (j-1)^2 \log \left(\prod_l N_l \right) + C_3 m + \log \left(\prod_l N_l \right) \right\} \\
 & \leq \sum_{j \geq K} C_1 \exp \left\{ -\frac{\widetilde{C}_2 (j-1)^2 \log \left(\prod_l N_l \right)}{j} + C_3 m + \log \left(\prod_l N_l \right) \right\} \rightarrow 0
 \end{aligned}$$

as $K \rightarrow \infty$, where the equality holds as $a_{NT} = \log(\prod_l N_l)/\sqrt{T}$, the last inequality holds when $j \geq 1$, and the limit holds by the same calculation in (A.44). Therefore, by substituting the tail bounds for the first term in (A.51), we obtain that $P \left\{ a_{NT}^{-1} \sup_{i_l} \sqrt{d_{i_l}(\widehat{\Theta}_{\cdot i_l}, \widehat{\Theta}_{\cdot i_l}^\sigma)} \geq K \right\} \rightarrow 0$ as $K \rightarrow \infty$.

This implies that $\sup_{i_l} d_{i_l}(\widehat{\Theta}_{\cdot i_l}, \widehat{\Theta}_{\cdot i_l}^\sigma) = O_p(a_{NT}^2)$. By using Lemma 23, we have that

$$\begin{aligned}
 & \sup_{i_l} \left\{ \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-(l)}}^-, \widehat{\mathcal{G}}_{-l}) - \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^-, \mathcal{G}_{-l}^0) \right\} \\
 & = O_p(T^{-1}(\sum_l \log N_l)^2).
 \end{aligned}$$

Therefore, together with (A.50), we have

$$\begin{aligned}
 & \sup_{i_l} \left\{ \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\widehat{\boldsymbol{\xi}}_{\widehat{g}_{i_l}}^{(l)}; \widehat{\boldsymbol{\xi}}_{g^{-(l)}}^-, \widehat{\mathcal{G}}_{-l}) - \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)0}; \boldsymbol{\xi}_{g^{-(l)}}^-, \mathcal{G}_{-l}^0) \right\} \\
 & = O_p(T^{-1}(\sum_l \log N_l)^2). \tag{A.55}
 \end{aligned}$$

By Lemma 23, we get

$$\begin{aligned}
 & \sup_{i_l} \left\{ \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{i_{-l}} |\widehat{\alpha}_{\widehat{g}_{i_1}^{(1)} \dots \widehat{g}_{i_q}^{(q)}} - \alpha_{g_{i_1}^{(1)0} \dots g_{i_q}^{(q)0}}|^2 \right\} \\
 & = O_p\{\tau_{\min}^{-1} T^{-1}(\sum_l \log N_l)^2\} = O_p\{T^{-1}(\sum_l \log N_l)^2\}.
 \end{aligned}$$

I.4 Proof of Theorem 5

The Figure A.6 visualizes the strong group memberships consistency when true group number $G_{1,0} = 3$ and estimated group number $\widehat{G}_1 = 4$.

To prove the result, we use Lemma 25 by setting $\boldsymbol{\xi} = \widehat{\boldsymbol{\xi}}$. First we note that $\widehat{\Theta}$ satisfies (A.88) by Theorem 2 and the condition for c_{gap} . In the following we first show that $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(l)}$

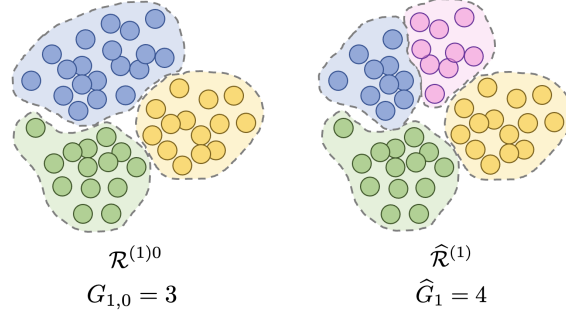


Figure A.6: The group memberships when $G_{1,0} = 3$ and $\widehat{G}_1 = 4$. The units belonging to the blue group are split into two sub groups (marked as blue and pink).

with $\eta = \tau_{\min} c_{\text{gap}} (c_{\pi})^{q-1} / \{8(\tau_{\min} + \tau_{\max})\}$ with probability tending to one. Then the conclusions of Lemma 25 hold. Second, we use the conclusion to prove the result. We remark that the notations in the following proof are borrowed in the statement in Lemma 25.

Step 1. Proof of $\widehat{\xi} \in \mathcal{N}_{\eta}^{(l)}$ with probability tending to one.

First by Lemma 24, we have

$$\begin{aligned} & \max_{g^{(l)} \in [G_{l,0}]} \min_{\tilde{g}^{(l)} \in [G_l]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \tilde{g}_{i-l}^{-(l)}}(\boldsymbol{\xi}) - \alpha_{g^{(l)} g_{i-l}^{-(l)0}}^0|^2 \right) \\ & = O_p(T^{-1} (\sum_l \log N_l)^2). \end{aligned}$$

Next we show

$$\begin{aligned} & \max_{\tilde{g}^{(l)} \in [G_l]} \min_{g^{(l)} \in [G_{l,0}]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \tilde{g}_{i-l}^{-(l)}}(\boldsymbol{\xi}) - \alpha_{g^{(l)} g_{i-l}^{-(l)0}}^0|^2 \right) \\ & = O_p(T^{-1} (\sum_l \log N_l)^2). \end{aligned}$$

Define the set $\mathcal{O}_{\tilde{g}^{(l)}} = \{i_l \in [N_l] : \widehat{g}_{i_l}^{(l)} = \tilde{g}^{(l)}\}$, where $\widehat{g}_{i_l}^{(l)}$ is obtained by (14) by using $\widehat{\boldsymbol{\xi}}$. Therefore for all $i_l \in \mathcal{O}_{\tilde{g}^{(l)}}$, it holds

$$\begin{aligned} & \min_{g^{(l)} \in [G_{l,0}]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \tilde{g}_{i-l}^{-(l)}}(\boldsymbol{\xi}) - \alpha_{g^{(l)} g_{i-l}^{-(l)0}}^0|^2 \right) \\ & \leq \|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \tilde{g}_{i-l}^{-(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i-l}^{-(l)0}}^0|^2 \\ & = \|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \tilde{g}_{i-l}^{-(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i-l}^{-(l)0}}^0|^2. \end{aligned}$$

Then it yields

$$\begin{aligned}
 & \max_{\tilde{g}^{(l)} \in [G_l]} \min_{g^{(l)} \in [G_{l,0}]} \left(\|\widehat{\boldsymbol{\theta}}_{\tilde{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}^{(l)} \widehat{g}_{i-l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g^{(l)} g_{i-l}^{(l)0}}^0|^2 \right) \\
 & \leq \max_{\tilde{g}^{(l)} \in [G_l]} \sup_{i_l \in \mathcal{G}_{\tilde{g}^{(l)}}} \left\{ \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\widehat{g}_{i_l}^{(l)} \widehat{g}_{i-l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i-l}^{(l)0}}^0|^2 \right\} \\
 & = \max_{i_l \in [N_l]} \left\{ \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\widehat{g}_{i_l}^{(l)} \widehat{g}_{i-l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i-l}^{(l)0}}^0|^2 \right\} \\
 & = O_p \left\{ T^{-1} \left(\sum_l \log N_l \right)^2 \right\},
 \end{aligned}$$

where the last equation is obtained by Proposition 4. By the condition that $c_{\text{gap}} \gg T^{-1}(\sum_l \log N_l)^2$, and the definition of $\mathcal{N}_\eta^{(l)}$, we can conclude $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(l)}$ with probability tending to one.

Step 2. Proof of conclusion of Theorem 5.

Since $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(l)}$, we have that the two conclusions in Lemma 25 holds for $\widehat{\boldsymbol{\xi}}$. Then it should be equivalent to prove that for $i_{l1}, i_{l2} \in \widehat{\mathcal{R}}_{\tilde{g}^{(l)}}^{(l)}$ for some $\tilde{g}^{(l)} \in [G_l]$, it holds $i_{l1}, i_{l2} \in \mathcal{R}_{g^{(l)}}^{(l)0}$ for some $g^{(l)0} \in [G_{l,0}]$. Suppose $i_{l1}, i_{l2} \in \widehat{\mathcal{R}}_{\tilde{g}^{(l)}}^{(l)}$ for some $\tilde{g}^{(l)} \in [G_l]$, then $\widehat{g}_{i_{l1}}^{(l)} = \widehat{g}_{i_{l2}}^{(l)} = \tilde{g}^{(l)}$ in this case. By conclusion (ii) in Lemma 25, we have $\tilde{g}^{(l)} \in \mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, g_{i_{l1}}^{(l)0}, \mathcal{G}_2^0)$ and $\tilde{g}^{(l)} \in \mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, g_{i_{l2}}^{(l)0}, \mathcal{G}_2^0)$ with probability tending to 1. Next, by conclusion (i) in Lemma 25, it holds $g_{i_{l1}}^{(l)0} = g_{i_{l2}}^{(l)0}$ since $\mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, \cdot, \mathcal{G}_2^0)$ is a partition of $[G_l]$. By defining $g^{(l)} \stackrel{\text{def}}{=} g_{i_{l1}}^{(l)0} = g_{i_{l2}}^{(l)0} \in [G_{l,0}]$, we have $i_{l1}, i_{l2} \in \mathcal{R}_{g^{(l)}}^{(l)0}$. Then the conclusion we finish the proof.

I.5 Proof of Corollary 6

By the Step 1 in the proof of Theorem 5, we have $\widehat{\boldsymbol{\xi}} \in \mathcal{N}_\eta^{(l)}$ with probability tending to one. For $G_l = G_{l,0}$, since $\{\mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, g^{(l)}, \mathcal{G}_2^0), g^{(l)} \in [G_{l,0}]\}$ is a partition of $[G_l]$, then we can conclude that $\mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, g^{(l)}, \mathcal{G}_2^0)$ contains only one element in $[G_l]$. Consequently, we can define a permutation by $\mathcal{A}_\eta^{(l)}(\widehat{\boldsymbol{\xi}}, \cdot, \mathcal{G}_2^0) : [G_{l,0}] \rightarrow [G_l]$. By Theorem 5, this implies that $\lim_{\min\{N_1, \dots, N_q, T\} \rightarrow \infty} P(\widehat{\mathcal{G}}_l = \mathcal{G}_l^0) \rightarrow 1$. Next, given that the probability of $\{\widehat{\mathcal{G}}_1 = \mathcal{G}_1^0, \dots, \widehat{\mathcal{G}}_q = \mathcal{G}_q^0\}$ goes to 1, (22) can be directly obtained.

I.6 Proof of Theorem 7

By Corollary 6, we have $\widehat{\boldsymbol{\xi}} = \widehat{\boldsymbol{\xi}}^{\text{or}}$ with probability tending to 1. Then it suffices to treat $\widehat{\boldsymbol{\xi}}$ as $\widehat{\boldsymbol{\xi}}^{\text{or}}$ equivalently. Recall that $\mathcal{E}_t^{(\mathcal{R}_{g^{(1)}}^{(1)}, \dots, \mathcal{R}_{g^{(q)}}^{(q)})} \in \mathbb{R}^{N_{1g^{(1)}} \times \dots \times N_{qg^{(q)}}}$ is defined as the subset of the tensor \mathcal{E}_t with each dimension selected by the index sets $\mathcal{R}_{g^{(l)}}^{(l)}$. By the vectorization of tensor defined in the main text, denote $\mathbb{Y}_{g^{(1)} \dots g^{(q)}, t} = \text{vec}(\mathcal{Y}_t^{(\mathcal{R}_{g^{(1)}}^{(1)}, \dots, \mathcal{R}_{g^{(q)}}^{(q)})})$ and $\mathbb{E}_{g^{(1)} \dots g^{(q)}, t} =$

$\text{vec}(\mathcal{E}_t^{(\mathcal{R}_{g^{(1)}}^{(1)}), \dots, \mathcal{R}_{g^{(q)}}^{(q)})}$, then we have

$$n^{q/2}T^{1/2}(\widehat{\boldsymbol{\xi}} - \boldsymbol{\xi}^0) = n^{q/2}T^{1/2}\mathbf{M}^{-1}\boldsymbol{\delta} = (n^{-q}T^{-1}\mathbf{M})^{-1}(n^{-q/2}T^{-1/2}\boldsymbol{\delta}),$$

where $\boldsymbol{\delta} = (\boldsymbol{\delta}^{(1)\top}, \dots, \boldsymbol{\delta}^{(q)\top}, \boldsymbol{\delta}^\alpha)^\top$ and

$$\begin{aligned} \boldsymbol{\delta}_{g^{(l)}}^{(l)} &= \sum_{t, g^{-(l)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{E}_{g^{(1)} \dots g^{(q)}, t}, & \boldsymbol{\delta}^{(l)} &= (\boldsymbol{\delta}_{g^{(l)}}^{(l)\top} : g^{(l)} \in [G_l])^\top \in \mathbb{R}^{G_l(p_l+1)}, \\ \boldsymbol{\delta}_{g^{(1)}, \dots, g^{(q)}}^\alpha &= \sum_t \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^\top \mathbb{E}_{g^{(1)} \dots g^{(q)}, t}, & \boldsymbol{\delta}^\alpha &= (\boldsymbol{\delta}_{g^{(1)}, \dots, g^{(q)}}^\alpha : \mathcal{I}_{g^{(1)}, \dots, g^{(q)}} \in \left[\prod_l G_l \right])^\top. \end{aligned}$$

Here

$$\begin{aligned} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} &= \left(\text{vec} \{ (\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)})_{g^{(1)} \dots g^{(q)}} \}, \right. \\ &\quad \left. \mathbf{1}_{N_{1g^{(1)}}} \otimes \dots \otimes (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}}^{(l)})} \otimes \dots \otimes \mathbf{1}_{N_{gg^{(q)}}} \right) \in \mathbb{R}^{(N_{1g^{(1)}} \dots N_{gg^{(q)}}) \times (p_l+1)} \end{aligned}$$

is defined in (10). Recall that $\mathbf{M}_{nT} = n^{-q}T^{-1}\mathbf{M}$ and then we have

$$\begin{aligned} n^{q/2}T^{1/2}\boldsymbol{\eta}^\top(\widehat{\boldsymbol{\xi}} - \boldsymbol{\xi}^0) &= \boldsymbol{\eta}^\top \mathbf{M}_{nT}^{-1} (n^{-q/2}T^{-1/2}\boldsymbol{\delta}) \\ &= n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta} + n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top \{ \mathbf{M}_{nT}^{-1} - (\mathbf{M}_{nT}^0)^{-1} \} \boldsymbol{\delta} \\ &= n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta} + n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top \mathbf{M}_{nT}^{-1} (\mathbf{M}_{nT}^0 - \mathbf{M}_{nT}) (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta}. \end{aligned}$$

Since we have $n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top \mathbf{M}_{nT}^{-1} (\mathbf{M}_{nT}^0 - \mathbf{M}_{nT}) (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta} = o_p(1)$, then it suffices to show

$$n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta} \rightarrow_d N\{0, \sigma^2 \boldsymbol{\eta}^\top (\mathbf{M}^0)^{-1}\boldsymbol{\eta}\}. \quad (\text{A.56})$$

Next, we provide the proof of (A.56).

Define $\tilde{\boldsymbol{\eta}} \stackrel{\text{def}}{=} (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\eta} = (\tilde{\boldsymbol{\eta}}_1^{(1)\top}, \dots, \tilde{\boldsymbol{\eta}}_{G_1}^{(1)\top}, \dots, \tilde{\boldsymbol{\eta}}_1^{(q)\top}, \dots, \tilde{\boldsymbol{\eta}}_{G_q}^{(q)\top}, \tilde{\boldsymbol{\eta}}_{1 \dots 1}^\alpha, \dots, \tilde{\boldsymbol{\eta}}_{G_1 \dots G_q}^\alpha)^\top$, where $\tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} \in \mathbb{R}^{p_l+1}$ and $\tilde{\boldsymbol{\eta}}_{g^{(1)} \dots g^{(q)}}^\alpha \in \mathbb{R}$. Then we have

$$\begin{aligned} n^{-q/2}T^{-1/2}\boldsymbol{\eta}^\top (\mathbf{M}_{nT}^0)^{-1}\boldsymbol{\delta} &= \frac{1}{\sqrt{nqT}} \sum_l \sum_{g^{(l)}} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)\top} \sum_{t, g^{-(l)}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{E}_{g^{(1)} \dots g^{(q)}, t} \\ &\quad + \frac{1}{\sqrt{nqT}} \sum_l \sum_{g^{(l)}=1}^{G_l} \tilde{\boldsymbol{\eta}}_{g^{(1)} \dots g^{(q)}}^\alpha \sum_t \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^\top \mathbb{E}_{g^{(1)} \dots g^{(q)}, t} \\ &\stackrel{\text{def}}{=} \sum_t \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^\top \mathbb{E}_{g^{(1)} \dots g^{(q)}, t} \stackrel{\text{def}}{=} \sum_t \mathcal{Z}_t^\top \mathbb{E}_t. \end{aligned}$$

Define \mathcal{F}_t as the sigma field generated by $\{\varepsilon_{i_1 \dots i_q, s} : i_l \in [N_l], l \in [q], s \leq t\}$. We have $\{\mathcal{Z}_t^\top \mathbb{E}_t, \mathcal{F}_t\}$ as a martingale difference array. By the central limit theorem of Hall and Heyde (2014), it suffices to verify for any $u > 0$

$$\sum_{t=1}^T E \left\{ (\mathcal{Z}_t^\top \mathbb{E}_t)^2 I(|\mathcal{Z}_t^\top \mathbb{E}_t| > u) | \mathcal{F}_{t-1} \right\} \rightarrow_p 0 \quad (\text{A.57})$$

$$\sum_{t=1}^T E \left\{ (\mathcal{Z}_t^\top \mathbb{E}_t)^2 | \mathcal{F}_{t-1} \right\} \rightarrow_p \boldsymbol{\eta}^\top (\mathbf{M}^0)^{-1} \mathbf{M}_b^0 (\mathbf{M}^0)^{-1} \boldsymbol{\eta}. \quad (\text{A.58})$$

where $\mathbf{M}_b^0 \stackrel{\text{def}}{=} \lim_{\min(n,T) \rightarrow \infty} n^{-q} T^{-1} \text{var}(\boldsymbol{\delta}) = \sigma^2 \lim_{\min(n,T) \rightarrow \infty} \mathbf{M}_{nT}^0 = \sigma^2 \mathbf{M}^0$, and $(\mathbf{M}^0)^{-1} = (\mathbf{M}^0)^{-1} \mathbf{M}_b^0 (\mathbf{M}^0)^{-1}$ by the condition that $\varepsilon_{i_1 \dots i_q, t}$ is i.i.d over i_1, \dots, i_q, t with $\text{var}(\varepsilon_{i_1 \dots i_q, t}) = \sigma^2$ and it is also independent of $\{\mathcal{Y}_{s-1}, \mathbf{X}_s^{(1)}, \dots, \mathbf{X}_s^{(q)} : s \leq t\}$, where $\mathbf{X}_s^{(l)} = (\mathbf{x}_{i_l, s}^{(l)} : i_l \in [N_l])^\top$.

(1) Proof of (A.57)

We have

$$\sum_{t=1}^T E \left\{ (\mathcal{Z}_t^\top \mathbb{E}_t)^2 I(|\mathcal{Z}_t^\top \mathbb{E}_t| > u) | \mathcal{F}_{t-1} \right\} \leq u^{-2} \sum_t E \left\{ (\mathcal{Z}_t^\top \mathbb{E}_t)^4 | \mathcal{F}_{t-1} \right\} \leq cu^{-2} \sum_t (\mathcal{Z}_t^\top \mathcal{Z}_t)^2,$$

where c is the fourth moment of $\varepsilon_{i_1 \dots i_q, t}$, which is a finite constant by Assumption 3. Then it suffices to show $E\{\sum_t (\mathcal{Z}_t^\top \mathcal{Z}_t)^2\} \rightarrow 0$. Note that we have

$$\begin{aligned} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t} &= \sum_l \frac{1}{\sqrt{n^q T}} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} + \frac{1}{\sqrt{n^q T}} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)} \tilde{\boldsymbol{\eta}}_{g^{(1)} \dots g^{(q)}}^\alpha \\ &\stackrel{\text{def}}{=} \sum_l \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)} + \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^\alpha. \end{aligned}$$

Then it suffices to show

$$E \left\{ \sum_t \left(\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \right)^2 \right\} \rightarrow 0, \quad (\text{A.59})$$

$$E \left\{ \sum_t \left(\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\alpha\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^\alpha \right)^2 \right\} \rightarrow 0 \quad (\text{A.60})$$

due to Cauchy-Schwarz inequality. We first show (A.59) in Step (1.1), and then show (A.60) in Step (1.2).

(1.1) Proof of $E\{\sum_t (\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)})^2\} \rightarrow 0$

Write $\tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} = (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)}, \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top})^\top$. By the definition of $\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)}$ in the main text, we have

$$\begin{aligned} &\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)}(\mathcal{R}_{g^{(l)}}, \cdot) \\ &= \left(\sum_{i_l=1}^{N_l} Y_{i_1, \dots, i_l, \dots, i_q, (t-1)} (a_{s, i_l}^{(l)} / n_{l_s}) \right)_{i_1 \in [N_1], \dots, s \in \mathcal{R}_{g^{(l)}}, \dots, i_q \in [N_q]} \in \mathbb{R}^{N_1 \times \dots \times N_{l_g^{(l)}} \times N_q}. \end{aligned}$$

We have

$$E \left\{ \sum_t \left(\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \right)^2 \right\} \leq c \frac{1}{n^{2q} T} E \left\{ \left(\sum_{g^{(1)}, \dots, g^{(q)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \|\mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\dagger}\|^2 \right)^2 \right\}$$

$$\begin{aligned}
 & + c \frac{1}{n^{2qT}} E \left\{ \left(\sum_{g^{(l)}, \dots, g^{(q)}} \left(\prod_{m \neq l} N_{mg^{(m)}} \right) \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\dagger\dagger\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\dagger\dagger} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right)^2 \right\} \\
 & = c \frac{1}{n^{2qT}} E \left\{ \left(\sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{vec} \{ \mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \}^\top \text{vec} \{ \mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \} \right)^2 \right\} \\
 & + c \frac{1}{n^{2qT}} E \left\{ \left(\sum_{g^{(l)}} \left(\prod_{m \neq l} N_m \right) \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})^\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right)^2 \right\},
 \end{aligned}$$

where c is a constant. By Lemma 22 and Cauchy-Schwarz inequality, it holds

$$\begin{aligned}
 & E \left\{ \left(\sum_g (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{vec} \{ \mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \}^\top \text{vec} \{ \mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \} \right)^2 \right\} \\
 & \leq c_4 \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{tr} \left(\mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} (\mathbf{1}_{N_1}^\top \cdots \mathbf{1}_{N_{l-1}}^\top \mathbf{1}_{N_l} \mathbf{1}_{N_{l+1}}^\top \cdots \mathbf{1}_{N_q}^\top) \right. \\
 & \quad \left. (\mathbf{1}_{N_q} \cdots \mathbf{1}_{N_{l+1}} \mathbf{1}_{N_l}^\top \mathbf{1}_{N_{l-1}} \cdots \mathbf{1}_{N_1}) \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})^\top} \right)^2 \leq c_4 \left(\sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 c_2 n^q \right)^2 \leq c n^{2q}
 \end{aligned}$$

where $c_4 = E(Y_{i_1 \dots i_q, t-1}^4) < \infty$ and c is a constant. Here the last second inequality holds because $N_l \leq c_2 n$ and c_2 is a positive constant. Next, we have

$$\begin{aligned}
 & \frac{1}{n^{2qT}} E \left\{ \left(\sum_{g^{(l)}} \left(\prod_{m \neq l} N_m \right) \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})^\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right)^2 \right\} \\
 & = \frac{1}{n^{2qT}} \left(\prod_{m \neq l} N_m \right)^2 E \left(\sum_{g^{(l)}} \sum_{i_t \in \mathcal{R}_{g^{(l)}}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_t}^{(l)})^2 \right)^2 \leq \frac{c}{n^{2T}} E \left(\sum_{g^{(l)}} \sum_{i_t \in \mathcal{R}_{g^{(l)}}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_t}^{(l)})^2 \right)^2.
 \end{aligned}$$

By the Cauchy's inequality we have

$$E \left\{ (\tilde{\boldsymbol{\eta}}_{g_1^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{1t}}^{(l)})^2 (\tilde{\boldsymbol{\eta}}_{g_2^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{2t}}^{(l)})^2 \right\} \leq \left[E \left\{ (\tilde{\boldsymbol{\eta}}_{g_1^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{1t}}^{(l)})^4 \right\} E \left\{ (\tilde{\boldsymbol{\eta}}_{g_2^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{2t}}^{(l)})^4 \right\} \right]^{1/2}.$$

We also have

$$\begin{aligned}
 E \left\{ (\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_t}^{(l)})^4 \right\} & = \left\| \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right\|^4 E \left\{ \left(\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\dagger\top} \mathbf{x}_{i_t}^{(l)} \right)^4 \right\} \\
 & = 4 \left\| \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right\|^4 \int_0^\infty t^3 P \left\{ |\langle \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\dagger}, \mathbf{x}_{i_t}^{(l)} \rangle| > t \right\} dt \leq c_K \left\| \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right\|^4
 \end{aligned}$$

where $\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\dagger} = \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} / \left\| \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right\|$, and c_K is a constant related to K . Here the last inequality is obtained by noting that $\langle \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\dagger}, \mathbf{x}_{i_t}^{(l)} \rangle$ is a 1-Lipschitz convex function of $\mathbf{x}_{i_t}^{(l)}$

with mean zero and $\mathbf{x}_{i_1 t}^{(l)}$ is K -convex defined in Definition 1. Consequently we have

$$\begin{aligned} & \frac{c}{n^2 T} E \left(\sum_{g_1^{(l)}, g_2^{(l)}} \sum_{i_{11} \in \mathcal{R}_{g_1^{(l)}}^{(l)}} \sum_{i_{12} \in \mathcal{R}_{g_2^{(l)}}^{(l)}} (\tilde{\boldsymbol{\eta}}_{g_1^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{11} t}^{(l)})^2 (\tilde{\boldsymbol{\eta}}_{g_2^{(l)}, (-1)}^{(l)\top} \mathbf{x}_{i_{12} t}^{(l)})^2 \right) \\ & \leq \frac{C}{n^2 T} \left(\sum_{g_1^{(l)}, g_2^{(l)}} c_K \|\tilde{\boldsymbol{\eta}}_{g_1^{(l)}, (-1)}^{(l)}\|^2 \|\tilde{\boldsymbol{\eta}}_{g_2^{(l)}, (-1)}^{(l)}\|^2 n^2 \right) = O(T^{-1}) = o(1). \end{aligned}$$

Consequently, it yields $E \left\{ \sum_t (\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{(l)})^2 \right\} = O(T^{-1}) = o(1)$.

(1.2) Proof of $E \left\{ \sum_t (\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\alpha\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\alpha})^2 \right\} \rightarrow 0$

Next, we have

$$\begin{aligned} & E \left\{ \sum_t \left(\sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\alpha\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\alpha} \right)^2 \right\} \\ & = \frac{1}{n^{2q} T} E \left\{ \sum_l \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(1)} \dots g^{(q)}}^{\alpha})^2 \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^{\top} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)} \right\}^2 \\ & \leq \frac{1}{n^{2q} T} c \left\{ \sum_l \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(1)} \dots g^{(q)}}^{\alpha})^2 n^q \right\}^2 = O(T^{-1}) = o(1), \end{aligned}$$

where the inequality is obtained by Lemma 22.

(2) Proof of (A.58)

We have

$$\begin{aligned} & \sum_{t=1}^T E \left\{ (\mathcal{Z}_t^{\top} \mathbb{E}_t)^2 | \mathcal{F}_{t-1} \right\} = \sigma^2 \sum_t \mathcal{Z}_t^{\top} \mathcal{Z}_t = \sigma^2 \sum_t \sum_l \sum_{g^{(l)}} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t}^{\top} \mathcal{Z}_{g^{(1)} \dots g^{(q)}, t} \\ & = \frac{\sigma^2}{n^q T} \sum_t \sum_l \sum_{g^{(1)}, \dots, g^{(q)}} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} \\ & + \frac{\sigma^2}{n^q T} \sum_t \sum_{g^{(1)}, \dots, g^{(q)}} (\tilde{\boldsymbol{\eta}}_{g^{(1)}, \dots, g^{(q)}}^{\alpha})^2 \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}^{\top} \mathbb{Y}_{g^{(1)} \dots g^{(q)}, (t-1)}. \end{aligned}$$

In the following we prove that

$$\begin{aligned} & \left| \frac{1}{n^q T} \sum_t \sum_{g^{(1)}, \dots, g^{(q)}} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} \right. \\ & \left. - \lim_{n \rightarrow \infty} \frac{1}{n^q} \sum_{g^{(1)}, \dots, g^{(q)}} \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)\top} E(\mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1)} \dots g^{(q)}, t}^{(l)}) \tilde{\boldsymbol{\eta}}_{g^{(l)}}^{(l)} \right| \rightarrow_p 0, \end{aligned} \quad (\text{A.61})$$

$$\begin{aligned} & \left| \frac{1}{n^q T} \sum_t \sum_{g^{(1), \dots, g^{(q)}}} (\tilde{\boldsymbol{\eta}}_{g^{(1), \dots, g^{(q)}}}^\alpha)^2 \mathbb{Y}_{g^{(1) \dots g^{(q)}, (t-1)}^\top \mathbb{Y}_{g^{(1) \dots g^{(q)}, (t-1)}} \right. \\ & \left. - \lim_{n \rightarrow \infty} \frac{1}{n^q} \sum_{g^{(1), \dots, g^{(q)}}} (\tilde{\boldsymbol{\eta}}_{g^{(1), \dots, g^{(q)}}}^\alpha)^2 E(\mathbb{Y}_{g^{(1) \dots g^{(q)}, (t-1)}^\top \mathbb{Y}_{g^{(1) \dots g^{(q)}, (t-1)}}) \right| \rightarrow_p 0. \end{aligned} \quad (\text{A.62})$$

We first show (A.61) in the following. Note that by the step (1.1), we have

$$\begin{aligned} & \frac{1}{n^q T} \sum_t \sum_{g^{(1), \dots, g^{(q)}}} \tilde{\boldsymbol{\eta}}_{g^{(1), \dots, g^{(q)}}}^{(l)\top} \mathbb{X}_{g^{(1) \dots g^{(q)}, t}^{(l)\top} \mathbb{X}_{g^{(1) \dots g^{(q)}, t}^{(l)} \tilde{\boldsymbol{\eta}}_{g^{(1), \dots, g^{(q)}}}^{(l)} \\ & = \frac{1}{n^q T} \sum_t \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{vec}\{\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})}\}^\top \text{vec}\{\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})}\} \\ & + \frac{1}{n^q T} \sum_t \sum_{g^{(l)}} \left(\prod_{m \neq l} N_m \right) \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})^\top} (\mathbf{X}_t^{(l)})^{(\mathcal{R}_{g^{(l)}, \cdot})} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)}. \end{aligned} \quad (\text{A.63})$$

Note that

$$\begin{aligned} & \frac{1}{n^q T} \sum_t \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{vec}\{\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})}\}^\top \text{vec}\{\mathcal{Y}_{t-1} \times_l \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})}\} \\ & = \frac{1}{n^q T} \sum_t \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \left((\mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \otimes \mathbf{I}_{N_{l+1}} \dots \otimes \mathbf{I}_{N_q}) \mathbb{Y}_t \right)^\top \\ & \left((\mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \otimes \mathbf{I}_{N_{l+1}} \dots \otimes \mathbf{I}_{N_q}) \mathbb{Y}_t \right). \end{aligned}$$

Further note that

$$\begin{aligned} & \left| \frac{1}{n^q T} \sum_t \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \left((\mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \otimes \mathbf{I}_{N_{l+1}} \dots \otimes \mathbf{I}_{N_q}) \mathbb{Y}_t \right)^\top \right. \\ & \left((\mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \otimes \mathbf{I}_{N_{l+1}} \dots \otimes \mathbf{I}_{N_q}) \mathbb{Y}_t \right) \\ & \left. - \frac{1}{n^q} \sum_{g^{(l)}} (\tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)})^2 \text{tr} \left\{ (\mathbf{I}_{N_q} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})^\top} \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \otimes \mathbf{I}_{N_{l+1}} \otimes \dots \otimes \mathbf{I}_{N_q}) \boldsymbol{\Gamma} \right\} \right| \\ & = o_p(1) \end{aligned} \quad (\text{A.64})$$

by Lemma 21, where $\widetilde{\mathbf{W}}$ in Lemma 21 is set to be $(\prod_{m \neq l} N_m)^{-1/2} (\mathbf{I}_{N_1} \otimes \dots \otimes \mathbf{I}_{N_{l-1}} \otimes \tilde{\boldsymbol{\eta}}_{g^{(l)}, 1}^{(l)} \mathbf{W}^{(l)(\mathcal{R}_{g^{(l)}, \cdot})}, 1 \leq g^{(l)} \leq G_l)^\top \otimes \mathbf{I}_{N_{l+1}} \otimes \dots \otimes \mathbf{I}_{N_q} \in \mathbb{R}^{(\prod_l N_l) \times (\prod_l N_l)}$ and we can verify that $n^{-1} \mathbf{1}_{\prod_l N_l}^\top \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}} \mathbf{1}_{\prod_l N_l}$ is bounded under when $N_l = N_{lg^{(l)}} = O(n)$ for all $l \in [q]$.

Next define $\mathbf{x}_t^{(l)\dagger} \stackrel{\text{def}}{=} ((\mathbf{X}_t^{(l)(\mathcal{R}_1, \cdot)} \tilde{\boldsymbol{\eta}}_{1, (-1)}^{(l)\dagger})^\top, (\mathbf{X}_t^{(l)(\mathcal{R}_2, \cdot)} \tilde{\boldsymbol{\eta}}_{2, (-1)}^{(l)\dagger})^\top, \dots, (\mathbf{X}_t^{(l)(\mathcal{R}_{G_l}, \cdot)} \tilde{\boldsymbol{\eta}}_{G_l, (-1)}^{(l)\dagger})^\top)^\top$, where $\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\dagger} = \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} / \|\tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)}\|$. Then we have

$$\frac{1}{nT} \sum_t \left(\sum_{g^{(l)}} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)\top} \mathbf{X}_t^{(l)(\mathcal{R}_{g^{(l)}, \cdot})^\top} \mathbf{X}_t^{(l)(\mathcal{R}_{g^{(l)}, \cdot})} \tilde{\boldsymbol{\eta}}_{g^{(l)}, (-1)}^{(l)} \right) = \frac{1}{nT} \sum_t \mathbf{x}_t^{(l)\dagger\top} \mathbf{A}^{(l)} \mathbf{x}_t^{(l)\dagger} \quad (\text{A.65})$$

where $\mathbf{A}^{(l)} = \text{diag}\{\|\tilde{\boldsymbol{\eta}}_{g^{(l)},(-1)}^{(l)}\|^2 \mathbf{1}_{N_{l_g^{(l)}}} : g^{(l)} \in [G_l]\}$. Note that $\mathbf{x}_t^{(l)\dagger}$ satisfies K -convex concentration property defined in Definition 1. Then by using Lemma 27, we have

$$|(nT)^{-1} \sum_t \mathbf{x}_t^{(l)\dagger\top} \mathbf{A}^{(l)} \mathbf{x}_t^{(l)\dagger} - n^{-1} E(\mathbf{x}_t^{(l)\dagger\top} \mathbf{A}^{(l)} \mathbf{x}_t^{(l)\dagger})| = o_p(1).$$

Together with (A.64) and (A.63), and by taking the limit, we have proved (A.61).

Next, (A.62) can be similarly proved by using the same technique as (A.64) and Lemma 21. Together with the fact that $n^{-q/2} T^{-1/2} \boldsymbol{\eta}^\top \mathbf{M}_{nT}^{-1} (\mathbf{M}_{nT}^0 - \mathbf{M}_{nT}) (\mathbf{M}_{nT}^0)^{-1} \boldsymbol{\delta} = o_p(1)$, we reach the conclusion.

Appendix J. Technical Lemmas

Lemma 17. *Suppose Assumptions 1–5 hold. Then we have*

$$\begin{aligned} & P\left(\sup_{\|\boldsymbol{\Theta}_{i_1 \dots i_q}\|_{\max} < R} \left| \frac{1}{T} Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) - \frac{1}{T} Q_{i_1 \dots i_q}^*(\boldsymbol{\Theta}_{i_1 \dots i_q}) \right| > x\right) \\ & \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_2 m\right\}, \end{aligned} \quad (\text{A.66})$$

where $m = \sum_l (p_l + 1) + 1$, and c_1, c_2 are positive constants, In addition, we have

$$\begin{aligned} & \sup_{\|\boldsymbol{\Theta}\|_{\max} < R} \left| \frac{1}{(\prod_l N_l) T} \{Q(\boldsymbol{\Theta}) - Q^*(\boldsymbol{\Theta})\} \right| \\ & \leq \sup_{i_1, \dots, i_q} \sup_{\|\boldsymbol{\Theta}_{i_1 \dots i_q}\|_{\max} < R} \left| \frac{1}{T} \{Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) - \frac{1}{T} Q_{i_1 \dots i_q}^*(\boldsymbol{\Theta}_{i_1 \dots i_q})\} \right| \\ & = O_p\left(T^{-1/2} (m + \sum_l \log N_l)\right). \end{aligned} \quad (\text{A.67})$$

Proof Note that we have

$$\begin{aligned} \frac{1}{T} Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) &= \frac{1}{T} \sum_{t=1}^T (\varepsilon_{i_1 \dots i_q, t} + \mathcal{X}_{i_1 \dots i_q, t}^\top \boldsymbol{\Theta}_{i_1 \dots i_q}^0 - \mathcal{X}_{i_1 \dots i_q, t}^\top \boldsymbol{\Theta}_{i_1 \dots i_q})^2 \\ &= \frac{1}{T} \sum_{t=1}^T \left\{ \varepsilon_{i_1 \dots i_q, t}^2 + 2\varepsilon_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\boldsymbol{\Theta}_{i_1 \dots i_q}^0 - \boldsymbol{\Theta}_{i_1 \dots i_q}) \right. \\ & \quad \left. + (\boldsymbol{\Theta}_{i_1 \dots i_q} - \boldsymbol{\Theta}_{i_1 \dots i_q}^0)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\boldsymbol{\Theta}_{i_1 \dots i_q} - \boldsymbol{\Theta}_{i_1 \dots i_q}^0) \right\}. \end{aligned}$$

Recall that $\boldsymbol{\Sigma}_{i_1 \dots i_q} = E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top)$. It is sufficient to show

$$P\left(\left|\frac{1}{T} \sum_t \varepsilon_{i_1 \dots i_q, t}^2 - \sigma^2\right| > x/3\right) \leq 2 \exp\{-c_1 T \min(x^2, x)\} \quad (\text{A.68})$$

$$P\left(\sup_{\|\mathbf{a}\|_{\max} < 2R} \left|\frac{2}{T} \sum_t \varepsilon_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \mathbf{a}\right| > x/3\right) \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_3 m\right\}, \quad (\text{A.69})$$

$$\begin{aligned}
 & P\left(\sup_{\|\mathbf{a}\|_{\max} < 2R} \left| \frac{1}{T} \sum_t \mathbf{a}^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \mathbf{a} - \mathbf{a}^\top \boldsymbol{\Sigma}_{i_1 \dots i_q} \mathbf{a} \right| > x/3\right) \\
 & \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_3 m\right\}
 \end{aligned} \tag{A.70}$$

where $m = \sum_l (p_l + 1) + 1$ as defined in Theorem 2, and c_1, c_3 are positive constants. First by Lemma 27, (A.68) holds directly. We next prove (A.70) in two steps. First we establish the upper bound for a fixed \mathbf{a} , then we establish the upper bound for all \mathbf{a} satisfying $\|\mathbf{a}\|_{\max} < 2R$. This yields (A.66) and (A.67) holds subsequently.

Step 1. We first show

$$P\left(\left|\frac{1}{T} \sum_t \mathbf{a}^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \mathbf{a} - \mathbf{a}^\top \boldsymbol{\Sigma}_{i_1 \dots i_q} \mathbf{a}\right| > x/3\right) \leq c_0 \exp\left(-c_3 \min(Tx^2, T^{1/2}x)\right)$$

for any \mathbf{a} , where c_3 and c_0 are positive constants. For the sake of similarity, we only verify the concentration inequality for diagonal elements of $T^{-1} \sum_t \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top$. Note that the diagonal elements of $T^{-1} \sum_t \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top$ takes the form $T^{-1} \sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w}$ (with $\mathbf{w} \in \mathbb{R}^{\prod_l N_l}$ and $\|\mathbf{w}\|_1 = 1$), or $T^{-1} \sum_t (\mathbf{x}_{ij,t}^{(l)})^2$, where $\mathbf{x}_{ij,t}^{(l)}$ is the j th element of the p_l -length vector $\mathbf{x}_{i,t}^{(l)}$. For $T^{-1} \sum_t \mathbf{x}_{i,t}^{(l)\top} \mathbf{x}_{i,t}^{(l)}$. We use Lemma 27 and Assumption 4, and note that $c_0 \exp\{-c_3 T \min(x^2, x)\} \leq c_0 \exp\{-c_3 \min(Tx^2, T^{1/2}x)\}$. This yields the concentration inequality. For the $T^{-1} \sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w}$, we use Lemma 18 to obtain the result.

Step 2. Here we follow Lemma F.2 of Basu et al. (2015) to prove the result. First we define $\mathcal{M} = \{\mathbf{v} \in \mathbb{R}^m : \|\mathbf{v}\|_{\max} \leq 2R\}$, where $m = \sum_l (p_l + 1) + 1$. Let $\mathcal{A} = \{\mathbf{u}_1, \dots, \mathbf{u}_{|\mathcal{A}|}\}$ be a $R/5$ -net of \mathcal{M} , where $|\mathcal{A}| \leq 10^m$. Hence for any $\mathbf{v} \in \mathcal{M}$, there exists some $\mathbf{u}_i \in \mathcal{A}$ such that $\|\Delta \mathbf{v}\|_{\max} \leq R/5$, where $\Delta \mathbf{v} = \mathbf{v} - \mathbf{u}_i$. Define $\mathbf{M} = |T^{-1} \sum_t \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top - \boldsymbol{\Sigma}_{i_1 \dots i_q}|_e$.

$$\begin{aligned}
 \tau & \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{M}} |\mathbf{v}^\top \mathbf{M} \mathbf{v}| \leq \max_{1 \leq k \leq |\mathcal{A}|} |\mathbf{u}_k^\top \mathbf{M} \mathbf{u}_k| + 2 \sup_{\mathbf{v} \in \mathcal{M}} \max_{1 \leq k \leq |\mathcal{A}|} |\Delta \mathbf{v}^\top \mathbf{M} \mathbf{u}_k| + \sup_{\mathbf{v} \in \mathcal{M}} |\Delta \mathbf{v}^\top \mathbf{M} \Delta \mathbf{v}| \\
 & \leq 2 \max_{1 \leq k \leq |\mathcal{A}|} |\mathbf{u}_k^\top \mathbf{M} \mathbf{u}_k| + 2 \sup_{\mathbf{v} \in \mathcal{M}} |\Delta \mathbf{v}^\top \mathbf{M} \Delta \mathbf{v}| \leq 2 \max_{1 \leq k \leq |\mathcal{A}|} |\mathbf{u}_k^\top \mathbf{M} \mathbf{u}_k| + \tau/50,
 \end{aligned}$$

where the last step is because $10\Delta \mathbf{v} \in \mathcal{M}$. Thus $\tau \leq (100/49) \max_{1 \leq k \leq |\mathcal{A}|} |\mathbf{u}_k^\top \mathbf{M} \mathbf{u}_k|$. This leads to

$$\begin{aligned}
 & \sup_{\|\mathbf{a}\|_{\max} \leq 2R} P\left(\left|\frac{1}{T} \sum_t \mathbf{a}^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \mathbf{a} - \mathbf{a}^\top \boldsymbol{\Sigma}_{i_1 \dots i_q} \mathbf{a}\right| > x/3\right) \\
 & = P\left\{\sup_{\mathbf{v} \in \mathcal{M}} |\mathbf{v}^\top \mathbf{M} \mathbf{v}| > \frac{x}{3}\right\} \leq P\left\{(100/49) \max_{1 \leq k \leq |\mathcal{A}|} |\mathbf{u}_k^\top \mathbf{M} \mathbf{u}_k| > x/3\right\} \\
 & \leq c_0 |\mathcal{A}| \exp\left(-c_4 \min(T(49x/100)^2, T^{1/2}49x/100)\right) \leq \exp\left\{-c_1 \min(Tx^2, T^{1/2}x) + c_3 m\right\},
 \end{aligned}$$

where c_3, c_4 are positive constants. (A.69) can be proved similarly hence we skip the details. \blacksquare

Lemma 18. *Under Assumption 4, then we have*

$$P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w} - \mathbf{w}^\top \mathbf{\Gamma} \mathbf{w}\right| > u\right) \leq 2(q^2 + q + 1) \exp\left(-C_1 \min(Tu^2, T^{1/2}u)\right), \quad (\text{A.71})$$

where $\mathbf{w} \in \mathbb{R}^{\prod_l N_l}$ is element-wisely non-negative and $\|\mathbf{w}\|_1 = 1$ and $\mathbf{\Gamma} = \text{cov}(\mathbb{Y}_t)$, and C_1, C_2, C_3 are positive constants.

Proof By (A.6), we have

$$\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w} = \frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} + \frac{2}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w} + \frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w}.$$

Specifically, let $\mathbf{\Gamma}^c = \text{cov}(\mathbb{Y}_t^c)$ and $\mathbf{\Gamma}^e = \text{cov}(\mathbb{Y}_t^e)$. It suffices to show

$$P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} - \mathbf{w}^\top \mathbf{\Gamma}^c \mathbf{w}\right| > u/3\right) \leq 2q^2 \exp\{-c_1 \min(Tu^2, T^{1/2}u)\} \quad (\text{A.72})$$

$$P\left(\left|\frac{2}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w}\right| > u/3\right) \leq 2q \exp\{-c_2 \min(Tu^2, T^{1/2}u)\} \quad (\text{A.73})$$

$$P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w} - \mathbf{w}^\top \mathbf{\Gamma}^e \mathbf{w}\right| > u/3\right) \leq 2 \exp\left(-c_3 \min\left(Tu^2, T^{1/2}u\right)\right). \quad (\text{A.74})$$

where c_1, c_2, c_3 are positive constants. In summary we obtain the final conclusion.

1. Proof of (A.72)

Let $\mathbf{A} = \mathbf{w} \mathbf{w}^\top \in \mathbb{R}^{(\prod_l N_l) \times (\prod_l N_l)}$. By (A.6),

$$\begin{aligned} \frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} &= \frac{1}{T}\sum_{t=1}^T \mathbb{Y}_t^{c\top} \mathbf{A} \mathbb{Y}_t^c = \frac{1}{T}\sum_{t=1}^T \sum_{l_1=0}^t \sum_{s_2=0}^t \mathbf{c}_{l_1}^\top (\mathbf{B}_0^{t-l_1})^\top \mathbf{A} \mathbf{B}_0^{t-s_2} \mathbf{c}_{s_2} \\ &= \sum_{l_1=0}^T \sum_{s_2=0}^T \mathbf{c}_{l_1}^\top \left\{ \frac{1}{T} \sum_{t=\max\{1, l_1, s_2\}}^T (\mathbf{B}_0^{t-l_1})^\top \mathbf{A} \mathbf{B}_0^{t-s_2} \right\} \mathbf{c}_{s_2} \stackrel{\text{def}}{=} \sum_{l_1=0}^T \sum_{s_2=0}^T \mathbf{c}_{l_1}^\top \mathbf{D}_{l_1 s_2} \mathbf{c}_{s_2}. \end{aligned} \quad (\text{A.75})$$

For any $l \geq 0$, we define $\mathbf{c} = (\mathbf{c}_0^\top, \mathbf{c}_1^\top, \dots, \mathbf{c}_T^\top)^\top$, and

$$\mathbf{D} = \begin{pmatrix} \mathbf{D}_{00} & \mathbf{D}_{01} & \cdots & \mathbf{D}_{0T} \\ \mathbf{D}_{10} & \mathbf{D}_{11} & \cdots & \mathbf{D}_{1T} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{D}_{T0} & \mathbf{D}_{T1} & \cdots & \mathbf{D}_{TT} \end{pmatrix}. \quad (\text{A.76})$$

Then we have $T^{-1}\sum_{t=1}^T \mathbb{Y}_t^{c\top} \mathbf{A} \mathbb{Y}_t^c = \mathbf{c}^\top \mathbf{D} \mathbf{c} = \sum_{l=1}^q \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)} + 2\sum_{l \neq m} \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)}$, where $\mathbf{c}^{(l)} = (\mathbf{c}_0^{(l)\top}, \dots, \mathbf{c}_T^{(l)\top})$. Obviously,

$$P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} - \mathbf{w}^\top \mathbf{\Gamma}^c \mathbf{w}\right| > u/3\right)$$

$$\begin{aligned}
 &= P\left(\left|\sum_{l=1}^q \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)} + 2 \sum_{l \neq m} \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)} - E\left(\sum_{l=1}^q \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)} + 2 \sum_{l \neq m} \mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)}\right)\right| > u/3\right) \\
 &\leq \sum_l P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)} - E(\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)})| > 2u/3q(q+1)\} \\
 &+ \sum_{l \neq m} P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)} - E(\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)})| > 2u/3q(q+1)\}.
 \end{aligned}$$

Recall that $\mathbf{c}_t^{(l)}$ is defined as $\text{vec}(\mathbf{1}_{N_1} \circ \dots \circ \boldsymbol{\beta}_{X_{l,t}}^0 \circ \dots \circ \mathbf{1}_{N_q}) \in \mathbb{R}^{\prod_l N_l}$ in the notations, and denote $\mathbf{c}^{(l)} = (\mathbf{c}_t^{(l)\top} : 0 \leq t \leq T)^\top \in \mathbb{R}^{(\prod_l N_l)(T+1)}$. Also denote the $C_u u = 2u/3q(q+1)$. Due to the bounded assumption for parameters in Assumption 1, we treat $\mathbf{c}_t^{(l)}$ as $\tilde{\mathbf{x}}_t^{(l)\eta}$ in Lemma 29, which indicates that

$$\begin{aligned}
 &P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)} - E(\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(l)})| > C_u u\} \\
 &= P\left\{\left|\sum_{l_1, s_2=0}^T \{\mathbf{c}_{l_1}^{(l)\top} \mathbf{D}_{l_1 s_2} \mathbf{c}_{s_2}^{(l)} - E(\mathbf{c}_{l_1}^{(l)\top} \mathbf{D}_{l_1 s_2} \mathbf{c}_{s_2}^{(l)})\}\right| > C_u u\right\} \\
 &\leq 2 \exp\left(-\frac{1}{C} \min\left(\frac{(C_u u)^2}{\|\mathcal{D}^{(l)}\|_F^2}, \frac{C_u u}{\|\mathcal{D}^{(l)}\|}\right)\right) \leq 2 \exp\{-c_1 \min(u^2 T, u T^{1/2})\},
 \end{aligned}$$

where $\mathcal{D}^{(l)}$ is given in Lemma 29 and the last step is due to Lemma 30. Similar treatment leads to

$$P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)} - E(\mathbf{c}^{(l)\top} \mathbf{D} \mathbf{c}^{(m)})| > C_u u\} \leq 2 \exp\{-c_1 \min(u^2 T, u T^{1/2})\}$$

Combining the above results, we get

$$P\left(\left|\frac{1}{T} \sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} - \mathbf{w}^\top \boldsymbol{\Gamma} \mathbf{w}\right| > u/3\right) \leq 2q^2 \exp\{-c_1 \min(Tu^2, T^{1/2}u)\},$$

where c_1 is a constant.

2. Proof of (A.73)

We can obtain that $\frac{1}{T} \sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w} = \sum_{l_1=0}^T \sum_{s_2=0}^T \mathbf{c}_{l_1}^\top \mathbf{D}_{l_1 s_2} \mathbb{E}_{s_2} = \mathbf{c}^\top \mathbf{D} \mathbb{E}$, where $\mathbb{E} = (\mathbb{E}_0^\top, \dots, \mathbb{E}_T^\top)^\top \in \mathbb{R}^{(\prod_l N_l)(T+1)}$. Obviously,

$$P\left(\left|\frac{1}{T} \sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w}\right| > u/6\right) = P\left(\left|\sum_l \mathbf{c}^{(l)\top} \mathbf{D} \mathbb{E}\right| > u/6\right) \leq \sum_l P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbb{E}| > u/(6q)\}.$$

Treating $\mathbf{c}_t^{(l)}$, which is defined as $\text{vec}(\mathbf{1}_{N_1} \circ \dots \circ \boldsymbol{\beta}_{X_{l,t}}^0 \circ \dots \circ \mathbf{1}_{N_q})$, as $\tilde{\mathbf{x}}_t^{(l)\eta}$ in Lemma 29, we get

$$\begin{aligned}
 &P\{|\mathbf{c}^{(l)\top} \mathbf{D} \mathbb{E}| > u/(6q)\} = P\left\{\left|\sum_{l_1, s_2=0}^T |\mathbf{c}_{l_1}^{(l)\top} \mathbf{D}_{l_1 s_2} \mathbb{E}_{s_2}|\right| > u/(6q)\right\} \\
 &\leq 2 \exp\left(-\frac{1}{C} \min\left(\frac{u^2}{\|\tilde{\mathcal{D}}^{(l)}\|_F^2}, \frac{u}{\|\tilde{\mathcal{D}}^{(l)}\|}\right)\right) \leq 2 \exp\{-c_2 \min(Tu^2, T^{1/2}u)\},
 \end{aligned}$$

where $\tilde{\mathcal{D}}^{(l)}$ is given in Lemma 29 and the last step is due to Lemma 30. Therefore, we get

$$P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w}\right| > u/6\right) \leq 2q \exp\{-c_2 \min(Tu^2, T^{1/2}u)\},$$

where c_2 is a positive constant.

3. Proof of (A.74)

Similar to the expression of (A.75), we can obtain that

$$\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w} = \sum_{l_1=0}^T \sum_{s_2=0}^T \mathbb{E}_{l_1}^\top \mathbf{D}_{l_1 s_2} \mathbb{E}_{s_2} = \mathbb{E}^\top \mathbf{D} \mathbb{E},$$

where $\mathbb{E} = (\mathbb{E}_0^\top, \dots, \mathbb{E}_T^\top)^\top \in \mathbb{R}^{\prod_l N_l}$. Since \mathbb{E} follows the K -convex concentration property, following Lemma 27, we have

$$\begin{aligned} P\left(\left|\frac{1}{T}\sum_t \mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w} - \mathbf{w}^\top \mathbf{\Gamma}^e \mathbf{w}\right| > u/3\right) &= P\left(\left|\mathbb{E}^\top \mathbf{D} \mathbb{E} - E(\mathbb{E}^\top \mathbf{D} \mathbb{E})\right| > u/3\right) \\ &\leq 2 \exp\left(-\frac{1}{C} \min\left(\frac{u^2/9}{K^4 \|\mathbf{D}\|_F^2}, \frac{u/3}{K^2 \|\mathbf{D}\|}\right)\right) \leq 2 \exp\left(-c_3 \min(Tu^2, T^{1/2}u)\right), \end{aligned}$$

where c_3 is a finite constant and the last step is due to (A.100) of Lemma 30. ■

Lemma 19. *Under Assumptions 1–5, denote $d_{i_1 \dots i_q} \stackrel{\text{def}}{=} d_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}, \Theta_{i_1 \dots i_q}^0)$. Recall that $S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) = Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}^0)$. we have*

$$\begin{aligned} P\left\{\sup_{\Theta_{i_1 \dots i_q}} T^{-1} \frac{|S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})|}{\sqrt{d_{i_1 \dots i_q}}} > x\right\} \\ \leq C_1 \exp\left\{-C_2 \min(Tx^2, \sqrt{T}x) + C_3 m\right\}, \end{aligned} \quad (\text{A.77})$$

where $m = \sum_l (p_l + 1) + 1$. Furthermore, we have

$$\max_l \sup_{i_l} \sup_{d_{i_1 \dots i_q} \leq \omega^2} T^{-1} |S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q})| = O_p\left\{\omega \left(\sum_l \log N_l\right) / \sqrt{T}\right\} \quad (\text{A.78})$$

$$\sup_{d(\Theta, \Theta^0) \leq \omega^2} \left(\prod_l N_l T\right)^{-1} |S(\Theta) - S^*(\Theta)| = O_p\left\{\omega \left(\sum_l \log N_l\right) / \sqrt{T}\right\}. \quad (\text{A.79})$$

Proof

1. Proof of (A.77).

We first write

$$\begin{aligned} &S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) - S_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q}) \\ &= \left[(\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)\right] \end{aligned}$$

$$- E\{(\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)\} \quad (\text{A.80})$$

$$+ 2\mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0) \varepsilon_{i_1 \dots i_q, t}. \quad (\text{A.81})$$

We use similar techniques in the proof of Lemma 17 to derive the concentration inequality. For example, we can apply (A.70) on term (A.80) by noticing that $\|\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0\| / \sqrt{d_{i_1 \dots i_q}} \leq cR$ for some constant c . Similarly, (A.81) can be shown by applying the proof for (A.69). Hence, we can obtain the conclusion under the same Assumptions of Lemma 17.

2. Proof of (A.78) and (A.79).

For simplicity, denote $d_{i_1 \dots i_q} \stackrel{\text{def}}{=} d_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}, \Theta_{i_1 \dots i_q}^0)$, and $S_{i_1 \dots i_q} \stackrel{\text{def}}{=} S_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q})$. Note that we have

$$\begin{aligned} & P\left\{ \sup_{d_{i_1 \dots i_q} \leq \omega^2} |S_{i_1 \dots i_q} - S_{i_1 \dots i_q}^*| > x \right\} \\ &= P\left\{ \sup_{d_{i_1 \dots i_q} \leq \omega^2} |S_{i_1 \dots i_q} - S_{i_1 \dots i_q}^*| / \omega > x / \omega \right\} \\ &\leq P\left\{ \sup_{d_{i_1 \dots i_q} \leq \omega^2} |S_{i_1 \dots i_q} - S_{i_1 \dots i_q}^*| / \sqrt{d_{i_1 \dots i_q}} > x / \omega \right\} \\ &\leq P\left\{ \sup_{\Theta_{i_1 \dots i_q}} |S_{i_1 \dots i_q} - S_{i_1 \dots i_q}^*| / \sqrt{d_{i_1 \dots i_q}} > x / \omega \right\} \\ &\leq C_1 \exp\left\{ -C_2 \min(Tx^2 / \omega^2, \sqrt{T}x / \omega) + C_3 m \right\}, \end{aligned}$$

where the last line uses (A.77). By taking the union bound, we have

$$\max_l \sup_{i_l} \sup_{d_{i_1 \dots i_q} \leq \omega^2} |S_{i_1 \dots i_q} - S_{i_1 \dots i_q}^*| = O_p\left\{ \omega \left(\sum_l \log N_l \right) / \sqrt{T} \right\}.$$

Using similar procedure, we can obtain (A.79) holds. ■

Lemma 20. *Under Assumptions 1–5, for vector $\Delta \in \mathbb{R}^{\sum_l (p_l + 1) + 1}$ in the parameter space, we have*

$$\begin{aligned} & P\left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \left(\prod_{m \neq l} N_m T \right)^{-1} \left| \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right\} \right| \geq x \right\} \\ &\leq C_1 \exp\left\{ -C_2 \min\left(Tx^2 / \omega^2, \sqrt{T}x / \omega\right) + C_3 m + \sum_l \log N_l \right\}, \quad (\text{A.82}) \end{aligned}$$

$$\begin{aligned} & P\left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \left(\prod_{m \neq l} N_m T \right)^{-1} \left| \sum_{m \neq l} \sum_{i_m} \sum_t \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta \varepsilon_{i_1 \dots i_q, t} \right| \geq x \right\} \\ &\leq C_1 \exp\left\{ -C_2 \min\left(Tx^2 / \omega^2, \sqrt{T}x / \omega\right) + C_3 m + \sum_l \log N_l \right\}, \quad (\text{A.83}) \end{aligned}$$

where $\mathcal{X}_{i_1 \dots i_q, t}$ is defined in equation (17), $m = \sum_l (p_l + 1) + 1$, and C_1, C_2, C_3 are positive constants.

Proof We first derive (A.82). Note that

$$\begin{aligned}
 & P \left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \frac{1}{\prod_{m \neq l} N_m T} \left| \sum_{m \neq l} \sum_{i_m} \sum_t \left\{ \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta \right. \right. \right. \\
 & \quad \left. \left. \left. - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right\} \right| \geq x \right\} \\
 & \leq P \left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} \right. \\
 & \quad \left. \frac{T^{-1} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|}{\|\Delta\|} \|\Delta\| \geq x \right\} \\
 & \leq P \left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \right. \\
 & \quad \left. \sqrt{\sum_{m \neq l} \sum_{i_m} \frac{T^{-2} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|^2}{\|\Delta\|^2}} \|\Delta\| \geq x \right\} \\
 & = P \left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \sum_{m \neq l} \sum_{i_m} \right. \\
 & \quad \left. \frac{T^{-2} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|^2}{\|\Delta\|^2} \|\Delta\|^2 \geq x^2 \right\} \\
 & \leq P \left\{ \sup_{i_l} \sup_{\|\Delta\|^2 \leq \omega^2} \sum_{m \neq l} \sum_{i_m} \right. \\
 & \quad \left. \frac{T^{-2} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|^2}{\|\Delta\|^2} \omega^2 \geq x^2 \right\} \\
 & \leq \sum_l \sum_{i_l=1}^{N_l} P \left\{ \sup_{\|\Delta\|^2 \leq \omega^2} \right. \\
 & \quad \left. \frac{T^{-2} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|^2}{\|\Delta\|^2} \omega^2 \geq x^2 \right\} \\
 & \leq \sum_l \sum_{i_l=1}^{N_l} P \left\{ \sup_{\Delta} \frac{T^{-1} \left| \sum_t \Delta^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top \Delta - [\Delta^\top E(\mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top) \Delta] \right|}{\|\Delta\|} \omega \geq x \right\} \\
 & \leq C_1 \exp \left\{ -C_2 \min(Tx^2/\omega^2, \sqrt{T}x/\omega) + C_3 m + \log \left(\prod_l N_l \right) \right\},
 \end{aligned}$$

where $m = \sum_l (p_l + 1) + 1$, the third line holds due to $\sum_i (a_i^2 b_i^2) \leq (\sum_i a_i^2)(\sum_i b_i^2)$. The last line is obtained by applying (A.70) and using that $\|\Delta\|_{\max}/\|\Delta\| \leq cR$. The proof of (A.83) can be finished in a similar scheme, while applying (A.69). \blacksquare

Lemma 21. *Under Assumption 4 and 5, let $\mathbf{M} = \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}} \in \mathbb{R}^{(\prod_l N_l) \times (\prod_l N_l)}$ be a symmetric matrix. Here $\widetilde{\mathbf{W}} \in \mathbb{R}^{m \times (\prod_l N_l)}$ is elementwisely positive and assume $n^{-1} \mathbf{1}_{\prod_l N_l}^\top \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}} \mathbf{1}_{\prod_l N_l} = O(1)$, where n is the order of N_l , $l \in [q]$. Then we have*

$$P\left(\left|\frac{1}{nT} \sum_t \mathbb{Y}_t^\top \mathbf{M} \mathbb{Y}_t - n^{-1} \text{tr}(\mathbf{M} \mathbf{\Gamma})\right| > u\right) \leq 2(q^2 + q + 1) \exp\left(-c \min(Tu^2, T^{1/2}u)\right),$$

and we have $n^{-1} \text{tr}(\mathbf{M} \mathbf{\Gamma}) \leq c$ for a positive constant c .

Proof By (A.6), we have $\frac{1}{T} \sum_t \mathbb{Y}_t^\top \mathbf{M} \mathbb{Y}_t = \frac{1}{T} \sum_t \mathbb{Y}_t^{c\top} \mathbf{M} \mathbb{Y}_t^c + \frac{2}{T} \sum_t \mathbb{Y}_t^{e\top} \mathbf{M} \mathbb{Y}_t^e + \frac{1}{T} \sum_t \mathbb{Y}_t^{e\top} \mathbf{M} \mathbb{Y}_t^e$. The proof follows the proof of Lemma 18. The difference is that we need to replace the \mathbf{A} matrix with \mathbf{M}/n . Since the procedure is the same we omit the details here.

Furthermore, we have $n^{-1} \text{tr}(\mathbf{M} \mathbf{\Gamma}) \leq n^{-1} \text{tr}(\mathbf{M} \mathbf{1}_{\prod_l N_l} \mathbf{1}_{\prod_l N_l}^\top \mathbf{c}_\Gamma) \leq c_\Gamma \mathbf{1}_{\prod_l N_l}^\top \mathbf{M} \mathbf{1}_{\prod_l N_l} / n = O(1)$, where the first inequality is obtained by Lemma 32. \blacksquare

Lemma 22. *Under Assumption 4, we have $\max_l \max_{i_l} E(Y_{i_1 \dots i_q, t}^4) \leq c$, where c is a positive constant.*

Proof Let $\mathbf{w} = \mathbf{e}_{i_q}^{(N_q)} \otimes \dots \otimes \mathbf{e}_{i_1}^{(N_1)} \in \mathbb{R}^{\prod_l N_l}$, where $\mathbf{e}_{i_l}^{(N_l)} \in \mathbb{R}^{N_l}$ is a vector whose i_l th element being equal to 1 while others being equal to 0. Then we have $Y_{i_1 \dots i_q, t} = \mathbf{w}^\top \mathbb{Y}_t$. By (A.6), we have

$$\begin{aligned} Y_{i_1 \dots i_q, t}^2 &= \mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w} = \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} + 2\mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{e\top} \mathbf{w} + \mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w} \\ &\leq 2\mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} + 2\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w}. \end{aligned}$$

Thus, $Y_{i_1 \dots i_q, t}^4 = (\mathbf{w}^\top \mathbb{Y}_t \mathbb{Y}_t^\top \mathbf{w})^2 \leq 4(\mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w})^2 + 4(\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^{e\top} \mathbf{w})^2$. It suffices to investigate the above two terms respectively.

1. Proof of $E\{(\mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w})^2\} < \infty$

First, let $\mathbf{A} = \mathbf{w} \mathbf{w}^\top \in \mathbb{R}^{\prod_l N_l \times \prod_l N_l}$, then by (A.6) and the definition in Lemma 30,

$$\begin{aligned} \mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^{c\top} \mathbf{w} &= \mathbb{Y}_t^{c\top} \mathbf{A} \mathbb{Y}_t^c = \sum_{s_1=0}^t \sum_{s_2=0}^t \mathbf{c}_{s_1}^\top (\mathbf{B}_0^{t-s_1})^\top \mathbf{A} \mathbf{B}_0^{t-s_2} \mathbf{c}_{s_2} \\ &= \sum_{s_1=0}^t \sum_{s_2=0}^t \mathbf{c}_{s_1}^\top \mathbf{L}_{t, s_1 s_2} \mathbf{c}_{s_2} = \sum_{s_1=0}^t \sum_{s_2=0}^t \left(\sum_l \mathbf{c}_{s_1}^{(l)\top} \right) \mathbf{L}_{t, s_1 s_2} \left(\sum_l \mathbf{c}_{s_2}^{(l)} \right) \\ &\leq 2 \sum_l \sum_{s_1=0}^t \sum_{s_2=0}^t \mathbf{c}_{s_1}^{(l)\top} \mathbf{L}_{t, s_1 s_2} \mathbf{c}_{s_2}^{(l)} = 2 \sum_l \sum_{s_1=0}^t \sum_{s_2=0}^t \beta_{X_l, s_1}^{(l)0\top} \mathcal{L}_{t, s_1 s_2}^{(l)} \beta_{X_l, s_2}^{(l)0}. \end{aligned} \quad (\text{A.84})$$

where $\mathcal{L}_{t, s_1 s_2}^{(l)} = (\mathbf{1}_{N_1}^\top \otimes \dots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \dots \otimes \mathbf{1}_{N_q}^\top) \mathbf{L}_{t, s_1 s_2} (\mathbf{1}_{N_1} \otimes \dots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \otimes \dots \otimes \mathbf{1}_{N_q}) \in \mathbb{R}^{N_l \times N_l}$ as defined in Lemma 30.

Further, let $\mathcal{L}_t^{(l)} = (\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} : 0 \leq s_1, s_2 \leq T)$, where $\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} = \mathbf{0}$ for $s_1 > t$ or $s_2 > t$, and $\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} = \mathcal{L}_{t, s_1 s_2}^{(l)}$ otherwise, as defined in Lemma 30. Let $\beta_{X_l}^{(l)0} = (\beta_{X_l, s}^{(l)0\top} : 0 \leq s \leq T)^\top$. Then we have $\sum_{s_1=0}^t \sum_{s_2=0}^t \beta_{X_l, s_1}^{(l)0\top} \mathcal{L}_{t, s_1 s_2}^{(l)} \beta_{X_l, s_2}^{(l)0} = \beta_{X_l}^{(l)0\top} \mathcal{L}_t^{(l)} \beta_{X_l}^{(l)0}$, which

implies that $(\mathbf{w}^\top \mathbb{Y}_t^c \mathbb{Y}_t^c \mathbf{w})^2 \leq 4 \sum_l (\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathcal{L}_t^{(l)} \boldsymbol{\beta}_{X_l}^{(l)0})^2$. Let $\mathcal{L}_t^{(l)} = \sum_{i_l} \lambda_{i_l} \mathbf{u}_{i_l} \mathbf{u}_{i_l}^\top$ be the eigen-decomposition of $\mathcal{L}_t^{(l)}$, where λ_{i_l} and \mathbf{u}_{i_l} are the i_l th eigenvalue and eigenvector respectively. One can verify that $\mathcal{L}_t^{(l)}$ is semi-definite, then we have all $\lambda_{i_l} \geq 0$. Then we have

$$(\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathcal{L}_t^{(l)} \boldsymbol{\beta}_{X_l}^{(l)0})^2 = \sum_{i_{l1}, i_{l2}} \lambda_{i_{l1}} \lambda_{i_{l2}} \left(\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathbf{u}_{i_{l1}} \mathbf{u}_{i_{l1}}^\top \boldsymbol{\beta}_{X_l}^{(l)0} \right) \left(\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathbf{u}_{i_{l2}} \mathbf{u}_{i_{l2}}^\top \boldsymbol{\beta}_{X_l}^{(l)0} \right).$$

Note that

$$\begin{aligned} & E \left\{ \left(\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathbf{u}_{i_{l1}} \mathbf{u}_{i_{l1}}^\top \boldsymbol{\beta}_{X_l}^{(l)0} \right) \left(\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathbf{u}_{i_{l2}} \mathbf{u}_{i_{l2}}^\top \boldsymbol{\beta}_{X_l}^{(l)0} \right) \right\} = E \left(\langle \mathbf{u}_{i_{l1}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle^2 \langle \mathbf{u}_{i_{l2}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle^2 \right) \\ & \leq (1/2) E \left(\langle \mathbf{u}_{i_{l1}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle^4 \right) + (1/2) E \left(\langle \mathbf{u}_{i_{l2}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle^4 \right) \\ & = 2 \int_0^\infty t^3 P \left\{ |\langle \mathbf{u}_{i_{l1}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle| > t \right\} dt + 2 \int_0^\infty t^3 P \left\{ |\langle \mathbf{u}_{i_{l2}}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle| > t \right\} dt \\ & \leq 8 \int_0^\infty t^3 \exp(-t^2/K^2) dt \stackrel{\text{def}}{=} c_K, \end{aligned}$$

where c_K is a constant related to K , and it is the constant defined in Definition 1. Here the last inequality is obtained by noting that $\langle \mathbf{u}_{i_l}, \boldsymbol{\beta}_{X_l}^{(l)0} \rangle$ is a 1-Lipschitz convex function of vector \mathbf{x} with mean zero and \mathbf{x} is K -convex defined in Definition 1. Therefore we have

$$E \left\{ (\boldsymbol{\beta}_{X_l}^{(l)0\top} \mathcal{L}_t^{(l)} \boldsymbol{\beta}_{X_l}^{(l)0})^2 \right\} \leq c_K \sum_{i_{l1}, i_{l2}} \lambda_{i_{l1}} \lambda_{i_{l2}} = c_K \text{tr}(\mathcal{L}_t^{(l)})^2 < \infty,$$

where the last inequality is obtained by (A.102) of Lemma 30.

2. Proof of $E\{(\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^e \mathbf{w})^2\} < \infty$

By the definition in Lemma 30, we have

$$\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^e \mathbf{w} = \mathbb{Y}_t^{e\top} \mathbf{A} \mathbb{Y}_t^e = \sum_{s_1=0}^t \sum_{s_2=0}^t \mathbb{E}_{s_1}^\top (\mathbb{B}_0^{t-s_1})^\top \mathbf{A} \mathbf{B}_0^{t-s_2} \mathbb{E}_{s_2} = \sum_{s_1=0}^t \sum_{s_2=0}^t \mathbb{E}_{s_1}^\top \mathbf{L}_{t, s_1 s_2} \mathbb{E}_{s_2}$$

Define $\tilde{\mathbf{L}}_t = (\tilde{\mathbf{L}}_{t, s_1 s_2} : 0 \leq s_1, s_2 \leq T)$, where $\tilde{\mathbf{L}}_{t, s_1 s_2} = \mathbf{0}$ for $s_1 > t$ or $s_2 > t$, and $\tilde{\mathbf{L}}_{t, s_1 s_2} = \mathbf{L}_{t, s_1 s_2}$ otherwise. Then we have $\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^e \mathbf{w} = \mathbb{E}^\top \tilde{\mathbf{L}}_t \mathbb{E}$, where $\mathbb{E} = (\mathbb{E}_0^\top, \dots, \mathbb{E}_T^\top)^\top$. Similar to the proof given in Step 1, we can verify $E\{(\mathbf{w}^\top \mathbb{Y}_t^e \mathbb{Y}_t^e \mathbf{w})^2\} = E\{(\mathbb{E}^\top \tilde{\mathbf{L}}_t \mathbb{E})^2\} \leq \text{ctr}(\tilde{\mathbf{L}}_t)^2 \leq \text{ctr}(\mathbf{L}_t)^2 < \infty$, where the second last inequality used the definition of \mathbf{L}_t in Lemma 30 and the last inequality is obtained by (A.102) of Lemma 30. \blacksquare

Lemma 23. *Under Assumption 2, it holds that*

$$\tau_{\min} \|\boldsymbol{\Theta}_{i_1 \dots i_q} - \boldsymbol{\Theta}_{i_1 \dots i_q}^0\|^2 \leq \frac{1}{T} \left\{ Q_{i_1 \dots i_q}^* (\boldsymbol{\Theta}_{i_1 \dots i_q}) - Q_{i_1 \dots i_q}^* (\boldsymbol{\Theta}_{i_1 \dots i_q}^0) \right\} \leq \tau_{\max} \|\boldsymbol{\Theta}_{i_1 \dots i_q} - \boldsymbol{\Theta}_{i_1 \dots i_q}^0\|^2 \quad (\text{A.85})$$

$$\tau_{\min} d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) \leq \frac{1}{(\prod_l N_l) T} \left\{ Q^*(\boldsymbol{\Theta}) - Q^*(\boldsymbol{\Theta}^0) \right\} \leq \tau_{\max} d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0). \quad (\text{A.86})$$

$$\begin{aligned}
 \tau_{\min} d_{i_l}(\Theta_{\cdot i_l}, \Theta_{\cdot i_l}^0) &\leq \frac{1}{(\prod_{m \neq l} N_l) T} \left\{ Q_{i_l}^*(\xi_{g_{i_l}}^{(l)}; \xi_{g^{-(l)}}^{-}, \mathcal{G}_{-l}) - Q_{i_l}^*(\xi_{g_{i_l}}^{(l)0}; \xi_{g^{-(l)0}}^{-}, \mathcal{G}_{-l}^0) \right\} \\
 &\leq \tau_{\max} d_{i_l}(\Theta_{\cdot i_l}, \Theta_{\cdot i_l}^0), \tag{A.87}
 \end{aligned}$$

where τ_{\min} is defined in Assumption 2, τ_{\max} is defined in notation Section A, and $d_{i_l}(\cdot, \cdot)$ is defined in (A.2).

Proof We have

$$\begin{aligned}
 Q(\Theta) &= \sum_{i_1, \dots, i_q} Q_{i_1 \dots i_q}(\Theta_{i_1 \dots i_q}) = \sum_l \sum_{i_l=1}^{N_l} \sum_{t=1}^T (\varepsilon_{i_1 \dots i_q, t} + \mathcal{X}_{i_1 \dots i_q, t}^\top \Theta_{i_1 \dots i_q}^0 - \mathcal{X}_{i_1 \dots i_q, t}^\top \Theta_{i_1 \dots i_q})^2 \\
 &= \sum_l \sum_{i_l=1}^{N_l} \sum_{t=1}^T \left\{ \varepsilon_{i_1 \dots i_q, t}^2 + 2\varepsilon_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q}^0 - \Theta_{i_1 \dots i_q}) \right. \\
 &\quad \left. + (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \mathcal{X}_{i_1 \dots i_q, t} \mathcal{X}_{i_1 \dots i_q, t}^\top (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0) \right\}.
 \end{aligned}$$

We have $\frac{1}{T} \left\{ Q_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q}) - Q_{i_1 \dots i_q}^*(\Theta_{i_1 \dots i_q}^0) \right\} = (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \Sigma_{i_1 \dots i_q} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)$.

By using Assumption 2 and Lemma 33, (A.85) can be obtained.

Next, it holds

$$\begin{aligned}
 &\frac{1}{(\prod_l N_l) T} \left\{ Q^*(\Theta) - Q^*(\Theta^0) \right\} \\
 &= \frac{1}{(\prod_l N_l) T} \sum_l \sum_{i_l=1}^{N_l} \sum_{t=1}^T (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \Sigma_{i_1 \dots i_q} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0) \\
 &= \frac{1}{\prod_l N_l} \sum_l \sum_{i_l=1}^{N_l} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \Sigma_{i_1 \dots i_q} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0).
 \end{aligned}$$

We have

$$\begin{aligned}
 &\frac{1}{\prod_l N_l} \sum_l \sum_{i_l=1}^{N_l} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \Sigma_{i_1 \dots i_q} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0) \\
 &\geq \frac{\tau_{\min}}{\prod_l N_l} \sum_l \sum_{i_l=1}^{N_l} \left\| \Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0 \right\|^2 = \tau_{\min} d(\Theta, \Theta^0), \\
 &\frac{1}{\prod_l N_l} \sum_l \sum_{i_l=1}^{N_l} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0)^\top \Sigma_{i_1 \dots i_q} (\Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0) \\
 &\leq \frac{\tau_{\max}}{\prod_l N_l} \sum_l \sum_{i_l=1}^{N_l} \left\| \Theta_{i_1 \dots i_q} - \Theta_{i_1 \dots i_q}^0 \right\|^2 = \tau_{\max} d(\Theta, \Theta^0)
 \end{aligned}$$

by using Assumption 2 and Lemma 33. This proves (A.86).

Similar arguments as that of the proof for (A.86) at a fixed i_l leads to (A.87). \blacksquare

Lemma 24. *Suppose $G_l \geq G_{l,0}$ for all $l \in [q]$, where $G_{l,0}$ is the true number of groups. In addition, assume Assumptions 1–5 hold. Define*

$$\sigma_l(g^{(l)}) = \operatorname{argmin}_{\widehat{g}^{(l)} \in [G_l]} \|\widehat{\boldsymbol{\theta}}_{\widehat{g}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\widehat{g}^{(l)} \widehat{g}_{i_m}^{-(l)}} - \alpha_{g^{(l)} g_{i_m}^{-(l)0}}^0 \right|^2.$$

Then we have

$$\begin{aligned} & \max_{g^{(l)} \in [G_{l,0}]} \left\{ \|\widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\sigma_l(g^{(l)}) \widehat{g}_{i_m}^{-(l)}} - \alpha_{g^{(l)} g_{i_m}^{-(l)0}}^0 \right|^2 \right\} \\ & \leq \max_{g^{(l)} \in [G_{l,0}]} \frac{N_l}{N_{lg^{(l)}}} d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}^0) = O_p\left(T^{-1} \left(\sum_l \log N_l \right)^2\right). \end{aligned}$$

Proof We have

$$\begin{aligned} & \|\widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\sigma_l(g^{(l)}) \widehat{g}_{i_m}^{-(l)}} - \alpha_{g^{(l)} g_{i_m}^{-(l)0}}^0 \right|^2 \\ & = \frac{1}{N_{lg^{(l)}}} \sum_{i_l} I(g_{i_l}^{(l)0} = g^{(l)}) \left\{ \|\widehat{\boldsymbol{\theta}}_{\sigma_l(g_{i_l}^{(l)0})}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 \right. \\ & \quad \left. + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\sigma_l(g_{i_l}^{(l)0}) \widehat{g}_{i_m}^{-(l)}} - \alpha_{g_{i_l}^{(l)0} g_{i_m}^{-(l)0}}^0 \right|^2 \right\} \\ & \leq \frac{1}{N_{lg^{(l)}}} \sum_{i_l} I(g_{i_l}^{(l)0} = g^{(l)}) \left\{ \|\widehat{\boldsymbol{\theta}}_{\widehat{g}_{i_l}^{(l)}}^{(l)} - \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\widehat{g}_{i_l}^{(l)} \widehat{g}_{i_m}^{-(l)}} - \alpha_{g_{i_l}^{(l)0} g_{i_m}^{-(l)0}}^0 \right|^2 \right\} \\ & \leq \frac{N_l}{N_{lg^{(l)}}} d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}^0). \end{aligned}$$

Then by taking $\max_{g^{(l)} \in [G_{l,0}]}$ of both sides, we have

$$\begin{aligned} & \max_{g^{(l)} \in [G_{l,0}]} \left(\|\widehat{\boldsymbol{\theta}}_{\sigma_l(g^{(l)})}^{(l)} - \boldsymbol{\theta}_{g^{(l)}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m} \left| \widehat{\alpha}_{\sigma_l(g^{(l)}) \widehat{g}_{i_m}^{-(l)}} - \alpha_{g^{(l)} g_{i_m}^{-(l)0}}^0 \right|^2 \right) \\ & \leq \frac{N_l}{\min_{g^{(l)}} N_{lg}} d(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}^0) = O_p\left(T^{-1} \left(\sum_l \log N_l \right)^2\right) \end{aligned}$$

by Theorem 2. ■

Lemma 25. *Under Assumptions 1–7, and suppose $G_l \geq G_{l,0}$ for all $l \in [q]$. For any $\boldsymbol{\xi}$, let $\widehat{g}_{i_l}^{(l)}(\boldsymbol{\xi})$ denote the membership obtained with (14). Define $\boldsymbol{\Theta} = (\boldsymbol{\Theta}_{i_1 \dots i_q} : i_l \in [N_l])$ with $\boldsymbol{\Theta}_{i_1 \dots i_q} = (\boldsymbol{\theta}_{\widehat{g}_{i_1}^{(1)}(\boldsymbol{\xi})}^{(1)\top}, \dots, \boldsymbol{\theta}_{\widehat{g}_{i_q}^{(q)}(\boldsymbol{\xi})}^{(q)\top}, \alpha_{\widehat{g}_{i_1}^{(1)}(\boldsymbol{\xi}) \dots \widehat{g}_{i_q}^{(q)}(\boldsymbol{\xi})}^{\top})^\top$. Assume we have*

$$d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) = O_p\left(T^{-1} \left(\sum_l \log N_l \right)^2\right) = o_p(c_{\text{gap}}), \quad (\text{A.88})$$

as $\min_l N_l \rightarrow \infty$. Let $\boldsymbol{\xi}^{(l)} = (\boldsymbol{\theta}^{(l)\top}, (\text{vec}(\boldsymbol{\alpha}_{\cdot g^{(l)}}))^\top : g^{(l)} \in [G_l])^\top = ((\boldsymbol{\theta}_{g^{(l)}}^{(l)} : g^{(l)} \in [G_l])^\top, (\text{vec}(\boldsymbol{\alpha}_{\cdot g^{(l)}}))^\top : g^{(l)} \in [G_l])^\top$. Define

$$\begin{aligned} \mathcal{M}^{(l)}(\boldsymbol{\xi}_{g^{(l)}}^{(l)}, \boldsymbol{\xi}_{g^{(l)0}}^{(l)0}; \mathcal{G}_{-l}, \mathcal{G}_{-l}^0) &= \|\boldsymbol{\theta}_{g^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)0}}^{(l)0}\|^2 \\ &+ \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m=1}^{N_m} |\alpha_{g_{i_l}^{(1)} \dots g_{i_{l-1}}^{(l-1)} g^{(l)} g_{i_{l+1}}^{(l+1)} \dots g_{i_q}^{(q)}} - \alpha_{g_{i_1}^0 \dots g_{i_{l-1}}^{(l-1)0} g^{(l)0} g_{i_{l+1}}^{(l+1)0} \dots g_{i_q}^{(q)0}}|^2 \\ &\stackrel{\text{def}}{=} \|\boldsymbol{\theta}_{g^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)0}}^{(l)0}\|^2 + \frac{1}{\prod_{m \neq l} N_l} \sum_{m \neq l} \sum_{i_m=1}^{N_m} |\alpha_{g^{(l)} g_{i_{-l}}^{-(l)}} - \alpha_{g^{(l)0} g_{i_{-l}}^{-(l)0}}|^2 \end{aligned}$$

with $\mathbf{i}_{-l} = (i_m : m \neq l)^\top$ and $g_{\mathbf{i}_{-l}}^{-(l)} = (g_{i_m}^{(m)} : m \neq l)^\top \in \mathbb{R}^{q-1}$. Further denote

$$\begin{aligned} d_S^{(l)}(\boldsymbol{\xi}^{(l)}, \boldsymbol{\xi}^{(l)0}; \mathcal{G}_{-l}, \mathcal{G}_{-l}^0) &= \max \left\{ \max_{g^{(l)0} \in [G_{l,0}]} \min_{g^{(l)} \in [G_l]} \left(\mathcal{M}^{(l)}(\boldsymbol{\xi}_{g^{(l)}}^{(l)}, \boldsymbol{\xi}_{g^{(l)0}}^{(l)0}; \mathcal{G}_{-l}, \mathcal{G}_{-l}^0) \right), \right. \\ &\quad \left. \max_{g^{(l)} \in [G_l]} \min_{g^{(l)0} \in [G_{l,0}]} \left(\mathcal{M}^{(l)}(\boldsymbol{\xi}_{g^{(l)}}^{(l)}, \boldsymbol{\xi}_{g^{(l)0}}^{(l)0}; \mathcal{G}_{-l}, \mathcal{G}_{-l}^0) \right) \right\}, \quad (\text{A.89}) \end{aligned}$$

where $\mathcal{G}_{-l} = \{\mathcal{G}_m : m \neq l\} = \{(g_{i_m}^{(m)} : 1 \leq i_m \leq N_m)^\top : m \neq l\}$. In addition, define $\mathcal{N}_\eta^{(l)} = \{\boldsymbol{\xi} : d_S^{(l)}(\boldsymbol{\xi}^{(l)}, \boldsymbol{\xi}^{(l)0}; \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}), \mathcal{G}_{-l}^0) < \eta\}$ given $\boldsymbol{\xi}$. Correspondingly, denote

$$\begin{aligned} \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g^{(l)0}, \mathcal{G}_{-l}^0) &= \left\{ g^{(l)} \in [G_l] : \|\boldsymbol{\theta}_{g^{(l)}}^{(l)} - \boldsymbol{\theta}_{g^{(l)0}}^{(l)0}\|^2 \right. \\ &\quad \left. + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g^{(l)} \widehat{g}_{\mathbf{i}_{-l}}^{-(l)}(\boldsymbol{\xi})} - \alpha_{g^{(l)0} g_{\mathbf{i}_{-l}}^{-(l)0}}| \leq \eta \right\}. \quad (\text{A.90}) \end{aligned}$$

Then the following conclusions hold:

- (i) For all $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(l)}$ with $\eta < (c_\pi)^{q-1} c_{\text{gap}}/4$, we have $\{\mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g^{(l)0}, \mathcal{G}_{-l}^0), g^{(l)0} \in [G_{l,0}]\}$ is a partition of $[G_l]$;
- (ii) Define the event $\boldsymbol{\Omega} = \{\widehat{g}_{i_l}^{(l)}(\boldsymbol{\xi}) \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0), \forall i_l \in [N_l]\}$, where $\boldsymbol{\xi}$ satisfying that $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(l)}$, and $\eta \leq \tau_{\min} c_{\text{gap}} (c_\pi)^{q-1} / \{8(\tau_{\min} + \tau_{\max})\}$. Here, τ_{\max} is given in the notation Section A. Then we have $P(\boldsymbol{\Omega}^c) \leq C \exp\left(-c_1 T^{1/2} c_{\text{gap}} + c_2 m + \sum_l \log N_l\right)$, where C, c_1, c_2 are positive constants.

Proof 1. Proof of (i)

By the definition of $\mathcal{N}_\eta^{(l)}$ and $\mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g^{(l)}, \mathcal{G}_{-l}^0)$, we have $\cup_{g^{(l)0} \in [G_{l,0}]} \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g^{(l)}, \mathcal{G}_{-l}^0) = [G_l]$. Then it remains to show that $\mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g^{(l)}, \mathcal{G}_{-l}^0)$ is a partition of $[G_l]$. That is $\mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_1^{(l)}, \mathcal{G}_{-l}^0) \cap \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_2^{(l)}, \mathcal{G}_{-l}^0) = \emptyset$ for any $g_1^{(l)} \neq g_2^{(l)}$.

We prove by contradiction. Suppose there exists $g_{12}^{(l)} \in [G_l]$ so that $g_{12}^{(l)} \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_1^{(l)}, \mathcal{G}_{-l}^0) \cap \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_2^{(l)}, \mathcal{G}_{-l}^0)$ for $g_1^{(l)}, g_2^{(l)} \in [G_{l,0}]$ and $g_1^{(l)} \neq g_2^{(l)}$. Denote $g^{-(l)} = (g^{(m)} : m \neq l)$, and correspondingly denote the event $\{g_{\mathbf{i}_{-l}}^{-(l)0} = g^{-(l)}\}$ as $\{g_{i_m}^{(m)0} = g^{(m)} : m \neq l\}$. In this case we have

$$c_{\text{gap}} \leq \left\{ \|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_2^{(l)}}^{(l)0}\|^2 + \max_{g^{-(l)} \in [G_{-l}^0]} |\alpha_{g_1^{(l)} g^{-(l)}}^0 - \alpha_{g_2^{(l)} g^{-(l)}}^0|^2 \right\}$$

$$\begin{aligned}
 &\leq 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \\
 &+ \max_{g^{-(l)} \in [G_{-l}^0]} \left(\frac{1}{\prod_{m \neq l} (\pi_{g^{(m)}, N_m}^{(m)} N_m)} \sum_{m \neq l} \sum_{i_m} \{ |\alpha_{g_1^{(l)} g^{-(l)}}^0 - \alpha_{g_2^{(l)} g^{-(l)}}^0|^2 I(g_{i_{-l}}^{-(l)0} = g^{-(l)}) \} \right) \\
 &= 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \\
 &+ \max_{g^{-(l)} \in [G_{-l}^0]} \left(\frac{1}{\prod_{m \neq l} (\pi_{g^{(m)}, N_m}^{(m)} N_m)} \sum_{m \neq l} \sum_{i_m} \{ |\alpha_{g_1^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_2^{(l)} g_{i_{-l}}^{-(l)0}}^0|^2 I(g_{i_{-l}}^{-(l)0} = g^{-(l)}) \} \right) \\
 &\leq 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \\
 &+ \max_{g^{-(l)} \in [G_{-l}^0]} \left(\frac{1}{\prod_{m \neq l} (\pi_{g^{(m)}, N_m}^{(m)} N_m)} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_1^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_2^{(l)} g_{i_{-l}}^{-(l)0}}^0|^2 \right) \\
 &= 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \\
 &+ \max_{g^{-(l)} \in [G_{-l}^0]} \left(\frac{1}{\prod_{m \neq l} (\pi_{g^{(m)}, N_m}^{(m)} N_m)} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_1^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_2^{(l)} g_{i_{-l}}^{-(l)0}}^0|^2 \right) \\
 &\leq 2\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + 2\|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \\
 &+ \frac{1}{\min_{g^{-(l)} \in [G_{-l}^0]} (\prod_{m \neq l} \pi_{g^{(m)}, N_m}^{(m)} N_m)} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_1^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_2^{(l)} g_{i_{-l}}^{-(l)0}}^0|^2 \\
 &\leq \frac{2}{\min_{g^{-(l)} \in [G_{-l}^0]} (\prod_{m \neq l} \pi_{g^{(m)}, N_m}^{(m)})} \times \left(\|\boldsymbol{\theta}_{g_1^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 + \|\boldsymbol{\theta}_{g_2^{(l)}}^{(l)0} - \boldsymbol{\theta}_{g_{12}^{(l)}}^{(l)}\|^2 \right) \\
 &+ \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_1^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_{12}^{(l)} \widehat{g}_{i_{-l}}^{-(l)}(\boldsymbol{\xi})}|^2 \\
 &+ \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_2^{(l)} g_{i_{-l}}^{-(l)0}}^0 - \alpha_{g_{12}^{(l)} \widehat{g}_{i_{-l}}^{-(l)}(\boldsymbol{\xi})}|^2 \\
 &\leq \frac{4\eta}{\min_{g^{-(l)} \in [G_{-l}^0]} (\prod_{m \neq l} \pi_{g^{(m)}, N_m}^{(m)})}.
 \end{aligned}$$

The last line contradicts the definition of η that $\eta < c_{\text{gap}}(c_\pi)^{q-1}/4$ as $N_l \rightarrow \infty$. Therefore, there does not exist a $g_{12}^{(l)} \in [G_l]$ such that $g_{12}^{(l)} \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_1^{(l)}, \mathcal{G}_{-l}^0) \cap \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_2^{(l)}, \mathcal{G}_{-l}^0)$ for any $g_1^{(l)}, g_2^{(l)} \in [G_{l,0}]$ and $g_1^{(l)} \neq g_2^{(l)}$, which suggests that $\mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_1^{(l)}, \mathcal{G}_{-l}^0) \cap \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_2^{(l)}, \mathcal{G}_{-l}^0) = \emptyset$ for any $g_1^{(l)} \neq g_2^{(l)}$. This completes the proof of part (i).

2. Proof of (ii)

For the notation simplicity, we use $\widehat{g}_{i_l}^{(l)}$ to replace $\widehat{g}_{i_l}^{(l)}(\boldsymbol{\xi})$, where $\boldsymbol{\xi}$ satisfies that $\boldsymbol{\xi} \in \mathcal{N}_\eta^{(l)}$ and also (A.88). By the definition of $\widehat{g}_{i_l}^{(l)}$, we have

$$I(\widehat{g}_{i_l}^{(l)} = g^{(l)}) \leq I\left(Q_{i_l}(\boldsymbol{\xi}_{g^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) < Q_{i_l}(\boldsymbol{\xi}_{g^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}))\right)$$

for any $\tilde{g}^{(l)} \neq g^{(l)}$. Therefore, for $\tilde{g}_{i_l}^{(l)} \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)$, we have

$$\begin{aligned} I\left(\tilde{g}_{i_l}^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) &= \sum_{g^{(l)=1}^{G_l}} I\left(g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) I(\tilde{g}_{i_l}^{(l)} = g^{(l)}) \\ &\leq \sum_{g^{(l)=1}^{G_l}} I\left(g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) I\left(Q_{i_l}(\boldsymbol{\xi}_{g^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) < Q_{i_l}(\boldsymbol{\xi}_{\tilde{g}_{i_l}^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}))\right) \\ &\stackrel{\text{def}}{=} \sum_{g^{(l)=1}^{G_l}} W_{i_l g^{(l)}}(\boldsymbol{\xi}). \end{aligned}$$

For all $g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)$ and $g^{(l)} \in [G_l]$, due to (i), there exists a $g_{j_l}^{(l)0} \neq g_{i_l}^{(l)0}$, such that $g^{(l)} \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{j_l}^{(l)0}, \mathcal{G}_{-l}^0)$. Then we have

$$\begin{aligned} &\left\| \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0} - \boldsymbol{\theta}_{g^{(l)}}^{(l)} \right\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g^{(l)} \tilde{g}_{i_l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i_l}^{(l)0}}^0|^2 \\ &\geq \frac{1}{2} \left\| \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0} - \boldsymbol{\theta}_{g_{j_l}^{(l)0}}^{(l)0} \right\|^2 + \frac{1}{2 \prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g_{j_l}^{(l)0} g_{i_l}^{(l)0}}^0 - \alpha_{g_{i_l}^{(l)0} g_{i_l}^{(l)0}}^0|^2 \\ &- \left\{ \left\| \boldsymbol{\theta}_{g_{j_l}^{(l)0}}^{(l)0} - \boldsymbol{\theta}_{g^{(l)}}^{(l)} \right\|^2 + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\alpha_{g^{(l)} \tilde{g}_{i_l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g_{j_l}^{(l)0} g_{i_l}^{(l)0}}^0|^2 \right\} \geq c_{\text{gap}}(c_\pi)^{q-1}/2 - \eta \end{aligned}$$

by Assumption 6 when $N_l \rightarrow \infty$. By Lemma 23, it holds for any $g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)$,

$$\frac{1}{(\prod_{m \neq l} N_m) T} \{Q_{i_l}^*(\boldsymbol{\xi}_{g^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) - Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)0}}, \mathcal{G}_{-l}^0)\} \geq \tau_{\min}(c_{\text{gap}}(c_\pi)^{q-1}/2 - \eta). \quad (\text{A.91})$$

On the other hand, for any $\tilde{g}_{i_l}^{(l)} \in \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)$, it holds

$$\begin{aligned} &\frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{\tilde{g}_{i_l}^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) - \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)0}}^{(l)0}; \boldsymbol{\xi}_{g^{-(l)0}}, \mathcal{G}_{-l}^0) \\ &\leq \tau_{\max} \left\{ \frac{1}{\prod_m N_m} \sum_{m \neq l} \sum_{i_m} \left\| \boldsymbol{\theta}_{\tilde{g}_{i_m}^{(m)}}^{(m)} - \boldsymbol{\theta}_{g_{i_m}^{(m)0}}^{(m)0} \right\|^2 + \left\| \boldsymbol{\theta}_{g_{i_l}^{(l)0}}^{(l)0} - \boldsymbol{\theta}_{\tilde{g}_{i_l}^{(l)}}^{(l)} \right\|^2 \right. \\ &\quad \left. + \frac{1}{\prod_{m \neq l} N_m} \sum_{m \neq l} \sum_{i_m} |\widehat{\alpha}_{\tilde{g}_{i_l}^{(l)} \tilde{g}_{i_l}^{(l)}}(\boldsymbol{\xi}) - \alpha_{g_{i_l}^{(l)0} g_{i_l}^{(l)0}}^0|^2 \right\} \\ &\leq \tau_{\max} \{d(\boldsymbol{\Theta}, \boldsymbol{\Theta}^0) + \eta\} \leq \tau_{\max} \{CT^{-1}(\sum_l \log N_l)^2 + \eta\}. \quad (\text{A.92}) \end{aligned}$$

in probability tending to 1 by (A.88), where C is a positive constant. Combining (A.91) and (A.92), we have

$$\frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{g^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) - \frac{1}{(\prod_{m \neq l} N_m) T} Q_{i_l}^*(\boldsymbol{\xi}_{\tilde{g}_{i_l}^{(l)}}^{(l)}; \boldsymbol{\xi}_{g^{-(l)}}, \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}))$$

$$\geq \tau_{\min}\{c_{\text{gap}}(c_\pi)^{q-1}/2 - \eta\} - \tau_{\max}\{\eta + CT^{-1}(\sum_l \log N_l)^2\} \stackrel{\text{def}}{=} \epsilon_\eta$$

when $N_l \rightarrow \infty$. Note that τ_{\min} and τ_{\max} are both bounded positive constants due to Assumption 2 and Lemma 33. This leads to

$$\begin{aligned} & W_{i_l g^{(l)}}(\mathbf{x}) \\ &= I\left(g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) I\left(Q_{i_l}(\boldsymbol{\xi}_{g^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) < Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}))\right) \\ &\leq I\left(g^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) \\ & I\left(\frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}(\boldsymbol{\xi}_{g_{i_l}^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) - \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}(\boldsymbol{\xi}_{g^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi}))\right) \\ &+ \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\boldsymbol{\xi}_{g^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) - \frac{1}{(\prod_{m \neq l} N_m)T} Q_{i_l}^*(\boldsymbol{\xi}_{g_{i_l}^{(l)}}; \boldsymbol{\xi}_{g^{-(l)}}^{-l}), \widehat{\mathcal{G}}_{-l}(\boldsymbol{\xi})) \geq \epsilon_\eta \\ &\leq 2I\left(\sup_{i_1, \dots, i_q} \sup_{\|\boldsymbol{\Theta}_{i_1, \dots, i_q}\|_{\max} < R} \left| \frac{1}{T} Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) - \frac{1}{T} Q_{i_1 \dots i_q}^*(\boldsymbol{\Theta}_{i_1 \dots i_q}) \right| \geq \epsilon_\eta/2\right). \end{aligned} \quad (\text{A.93})$$

Consequently, we have

$$\begin{aligned} & P\left\{\sup_{1 \leq i_l \leq N_l} I\left(\widehat{g}_{i_l}^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) > 0\right\} \leq \sum_{g^{(l)=1}^{G_l} P\left\{\sup_{1 \leq i_l \leq N_l} W_{i_l g^{(l)}}(\boldsymbol{\xi}) > 0\right\} \\ &\leq \sum_{g^{(l)=1}^{G_l} P\left\{\sup_{i_1, \dots, i_q} \sup_{\|\boldsymbol{\Theta}_{i_1, \dots, i_q}\|_{\max} < R} \left| \frac{1}{T} Q_{i_1 \dots i_q}(\boldsymbol{\Theta}_{i_1 \dots i_q}) - \frac{1}{T} Q_{i_1 \dots i_q}^*(\boldsymbol{\Theta}_{i_1 \dots i_q}) \right| \geq \epsilon_\eta/2\right\} \\ &\leq G_l \left(\prod_l N_l\right) \exp\left\{-c'_1 \min(T\epsilon_\eta^2, T^{1/2}\epsilon_\eta) + c_2 m\right\} \\ &= G_l \exp\left\{-c'_1 \min(T\epsilon_\eta^2, T^{1/2}\epsilon_\eta) + c_2 m + \sum_l \log N_l\right\} \end{aligned} \quad (\text{A.94})$$

by using (A.66) of Lemma 17. Finally, noting that $c_{\text{gap}} \gg T^{-1}(\sum_l \log N_l)^2$ and we have $\eta < \{\tau_{\min} c_{\text{gap}} (c_\pi)^{q-1}/4 - \tau_{\max} CT^{-1}(\sum_l \log N_l)^2\}/(\tau_{\min} + \tau_{\max})$, then we can obtain $\epsilon_\eta > \tau_{\min} c_{\text{gap}} (c_\pi)^{q-1}/4$ as $\min_l N_l, T \rightarrow \infty$. This results in

$$P\left\{\sup_{1 \leq i_l \leq N_l} I\left(\widehat{g}_{i_l}^{(l)} \notin \mathcal{A}_\eta^{(l)}(\boldsymbol{\xi}, g_{i_l}^{(l)0}, \mathcal{G}_{-l}^0)\right) > 0\right\} \leq C \exp\left\{-c_1 T^{1/2} c_{\text{gap}} + c_2 m + \sum_l \log N_l\right\}.$$

This finishes the proof of (ii). ■

Appendix K. General Technical Lemmas

Lemma 26. (HANSON-WRIGHT INEQUALITY) *Define the ϕ_2 -norm of variable X as*

$$\|X\|_{\phi_2} = \inf\left\{t > 0 : E \exp\left(\frac{X^2}{t^2}\right) \leq 2\right\}.$$

Let X_1, \dots, X_n be independent variables satisfying $E(X_i) = 0$, $E(X_i^2) = \sigma_i^2$, and $\|X_i\|_{\phi_2} \leq M < \infty$ and \mathbf{A} be a symmetric $n \times n$ matrix. Define $\mathbf{x} = (X_1, \dots, X_n)^\top \in \mathbb{R}^n$. For any $t > 0$, we have

$$P\left(|\mathbf{x}^\top \mathbf{A} \mathbf{x} - \text{tr}(\mathbf{A} \boldsymbol{\Sigma}_x)| > t\right) \leq 2 \exp\left\{-c \min\left(\frac{t^2}{M^4 \|\mathbf{A}\|_F^2}, \frac{t}{M^2 \sigma_1(\mathbf{A})}\right)\right\},$$

where $\boldsymbol{\Sigma}_x = \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$, $\sigma_1(\mathbf{A})$ is the maximum singular value of \mathbf{A} , $\|\mathbf{A}\|_F$ denotes the Frobenius norm of matrix \mathbf{A} , and $c > 0$ is a constant.

Proof The proof is given in Hanson and Wright (1971). ■

Lemma 27. *let \mathbf{x} be a mean zero random vector in \mathbb{R}^n satisfying K -convex concentration property according to Definition 1. Then for any $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $t > 0$, it holds*

$$P\left(|\mathbf{x}^\top \mathbf{A} \mathbf{x} - E(\mathbf{x}^\top \mathbf{A} \mathbf{x})| \geq t\right) \leq 2 \exp\left(-\frac{1}{C} \min\left(\frac{t^2}{K^4 \|\mathbf{A}\|_F^2}, \frac{t}{K^2 \|\mathbf{A}\|}\right)\right).$$

Proof The proof can be found in Theorem 2.5 of Adamczak (2015). ■

Lemma 28. *If $\mathbf{x} \in \mathbb{R}^n$ satisfies the K -convex concentration property according to Definition 1 with $E(\mathbf{x}) = \mathbf{0}$, then we have $\|\text{cov}(\mathbf{x})\| \leq 2K^2$.*

Proof For any unit vector $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\| = 1$, $\langle \mathbf{u}, \mathbf{x} \rangle$ is a 1-Lipschitz convex function of \mathbf{x} . Then we have

$$\mathbf{u}^\top \text{cov}(\mathbf{x}) \mathbf{u} = E\langle \mathbf{u}, \mathbf{x} \rangle^2 = 2 \int_0^\infty t P(|\langle \mathbf{u}, \mathbf{x} \rangle| > t) dt \leq 4 \int_0^\infty t \exp(-t^2/K^2) dt = 2K^2.$$

This implies $\|\text{cov}(\mathbf{x})\| \leq 2K^2$. ■

Lemma 29. *Let $\{\mathbf{D}_{s_1 s_2} \in \mathbb{R}^{(\prod_l N_l) \times (\prod_l N_l)} : s_1, s_2 \in \{0, \dots, T\}\}$ be a sequence of matrices. Define the following three terms, $\mathcal{D}_{s_1 s_2}^{(l)} = (\mathbf{1}_{N_1}^\top \otimes \dots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \dots \otimes \mathbf{1}_{N_q}^\top) \mathbf{D}_{s_1 s_2} (\mathbf{1}_{N_1} \otimes \dots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \otimes \dots \otimes \mathbf{1}_{N_q}) \in \mathbb{R}^{N_l \times N_l}$, $\mathcal{D}_{s_1 s_2}^{(lm)} = (\mathbf{1}_{N_1}^\top \otimes \dots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \dots \otimes \mathbf{1}_{N_q}^\top) \mathbf{D}_{s_1 s_2} (\mathbf{1}_{N_1} \otimes \dots \otimes \mathbf{1}_{N_{m-1}} \otimes \mathbf{I}_{N_m} \otimes \mathbf{1}_{N_{m+1}} \otimes \dots \otimes \mathbf{1}_{N_q}) \in \mathbb{R}^{N_l \times N_m}$, and $\tilde{\mathcal{D}}_{s_1 s_2}^{(l)} = (\mathbf{1}_{N_1}^\top \otimes \dots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \dots \otimes \mathbf{1}_{N_q}^\top) \mathbf{D}_{s_1 s_2} \in \mathbb{R}^{N_l \times (\prod_l N_l)}$. Correspondingly, let $\mathcal{D}^{(l)} = (\mathcal{D}_{s_1 s_2}^{(l)} : s_1, s_2 \in \{0, \dots, T\}) \in \mathbb{R}^{(N_l(T+1)) \times (N_l(T+1))}$, let $\mathcal{D}^{(lm)} = (\mathcal{D}_{s_1 s_2}^{(lm)} : s_1, s_2 \in \{0, \dots, T\}) \in \mathbb{R}^{(N_l(T+1)) \times (N_m(T+1))}$, and let $\tilde{\mathcal{D}}^{(l)} = (\tilde{\mathcal{D}}_{s_1 s_2}^{(l)} : s_1, s_2 \in \{0, \dots, T\}) \in \mathbb{R}^{N_l(T+1) \times (\prod_l N_l)(T+1)}$. Define $\tilde{\mathbf{x}}_t^{(l)\eta} = \mathbf{1}_{N_1} \circ \dots \circ \mathbf{x}_t^{(l)\eta} \circ \dots \circ \mathbf{1}_{N_q} \in \mathbb{R}^{\prod_l N_l}$, where $\mathbf{x}_t^{(l)\eta}$ is defined in Assumption 4. Further note that $\mathbf{x}^{(l)\eta} = (\mathbf{x}_t^{(l)\eta} : 0 \leq t \leq T) \in \mathbb{R}^{N_l(T+1)}$. Then we have*

$$P\left\{\left|\sum_{s_1, s_2=0}^T \tilde{\mathbf{x}}_t^{(l)\eta \top} \mathbf{D}_{s_1 s_2} \tilde{\mathbf{x}}_t^{(l)\eta} - \sum_{s_1, s_2=0}^T E(\tilde{\mathbf{x}}_t^{(l)\eta \top} \mathbf{D}_{s_1 s_2} \tilde{\mathbf{x}}_t^{(l)\eta})\right| \geq u\right\}$$

$$\leq 2 \exp \left(-\frac{1}{C} \min \left(\frac{u^2}{\|\mathcal{D}^{(l)}\|_F^2}, \frac{u}{\|\mathcal{D}^{(l)}\|} \right) \right), \quad (\text{A.95})$$

$$P \left\{ \left| \sum_{s_1, s_2=0}^T \tilde{\mathbf{x}}_{s_1}^{(l)\eta\top} \mathbf{D}_{s_1 s_2} \tilde{\mathbf{x}}_{s_2}^{(m)\eta} \right| > u \right\} \leq 2 \exp \left(-\frac{1}{C} \min \left(\frac{u^2}{\|\mathcal{D}^{(lm)}\|_F^2}, \frac{u}{\|\mathcal{D}^{(lm)}\|} \right) \right), \quad (\text{A.96})$$

$$P \left\{ \left| \sum_{s_1, s_2=0}^T \tilde{\mathbf{x}}_{s_1}^{(l)\eta\top} \mathbf{D}_{s_1 s_2} \mathbb{E}_{s_2} \right| > u \right\} \leq 2 \exp \left(-\frac{1}{C} \min \left(\frac{u^2}{\|\tilde{\mathcal{D}}^{(l)}\|_F^2}, \frac{u}{\|\tilde{\mathcal{D}}^{(l)}\|} \right) \right). \quad (\text{A.97})$$

Proof 1. Proof of (A.95)

Note that we have $\sum_{s_1, s_2=0}^T \tilde{\mathbf{x}}_{s_1}^{(l)\eta\top} \mathbf{D}_{s_1 s_2} \tilde{\mathbf{x}}_{s_2}^{(l)\eta} = \sum_{s_1, s_2=0}^T \mathbf{x}_{s_1}^{(l)\eta\top} \mathcal{D}_{s_1 s_2}^{(l)} \mathbf{x}_{s_2}^{(l)\eta} = \mathbf{x}^{(l)\eta\top} \mathcal{D}^{(l)} \mathbf{x}^{(l)\eta}$. Here $\mathbf{x}^{(l)\eta}$ satisfies the K' -convex concentration property defined in Definition 1 by Assumption 4 for a constant K' . Therefore by Lemma 27 the concentration inequality in (A.95) is directly obtained.

2. Proof of (A.96)

Note that $\sum_{s_1, s_2=0}^T \mathbf{x}_{s_1}^{(l)\eta\top} \mathbf{D}_{s_1 s_2} \mathbf{x}_{s_2}^{(m)\eta} = \sum_{s_1, s_2=0}^T \mathbf{x}_{s_1}^{(l)\eta\top} \mathcal{D}_{s_1 s_2}^{(lm)} \mathbf{x}_{s_2}^{(m)\eta} = \mathbf{x}^{(l)\eta\top} \mathcal{D}^{(lm)} \mathbf{x}^{(m)\eta}$. Define $\mathbf{h}^\eta = (\mathbf{x}^{(l)\eta\top}, \mathbf{x}^{(m)\eta\top})^\top$, and $\mathcal{D} = (\mathbf{0}, \mathcal{D}^{(lm)}; \mathcal{D}^{(lm)\top}, \mathbf{0})$. Then we have $\mathbf{x}^{(l)\eta\top} \mathcal{D}^{(lm)} \mathbf{x}^{(m)\eta} = \mathbf{h}^{\eta\top} \mathcal{D} \mathbf{h}^\eta / 2$. By Lemma 27 and Assumption 4 we obtain that

$$\begin{aligned} P \left(\left| \mathbf{x}^{(l)\eta\top} \mathcal{D}^{(lm)} \mathbf{x}^{(m)\eta} \right| > u \right) &\leq 2 \exp \left(-\frac{1}{C_1} \min \left(\frac{u^2}{\|\mathcal{D}\|_F^2}, \frac{u}{\|\mathcal{D}\|} \right) \right) \\ &\leq 2 \exp \left(-\frac{1}{C} \min \left(\frac{u^2}{\|\mathcal{D}^{(lm)}\|_F^2}, \frac{u}{\|\mathcal{D}^{(lm)}\|} \right) \right). \end{aligned}$$

3. Proof of (A.97)

Denote $\mathbb{E} = (\mathbb{E}_t^\top : 0 \leq t \leq T)^\top \in \mathbb{R}^{(\prod_l N_l)(T+1)}$. Note that $\sum_{s_1, s_2=0}^T \tilde{\mathbf{x}}_{s_1}^{(l)\eta\top} \mathbf{D}_{s_1 s_2} \mathbb{E}_{s_2} = \mathbf{x}^{(l)\eta\top} \tilde{\mathcal{D}}^{(l)} \mathbb{E}$, then we follow the proof of (A.96) and obtain the result. ■

Lemma 30. Let $\mathbf{D}_{s_1 s_2} = 1/T \sum_{t=\max\{1, s_1, s_2\}}^T (\mathbf{B}_0^{t-s_1})^\top \mathbf{w}^{(l)} \mathbf{w}^{(l)\top} \mathbf{B}_0^{t-s_2} \stackrel{\text{def}}{=} 1/T \sum_{t=\max\{1, s_1, s_2\}}^T \mathbf{L}_{t, s_1 s_2}^{(l)}$, where $\mathbf{w}^{(l)} = \mathbf{e}_{k_1}^{(N_1)} \otimes \cdots \otimes \mathbf{w}_{i_l}^{(l)} \otimes \cdots \otimes \mathbf{e}_{k_q}^{(N_q)} \in \mathbb{R}^{\prod_l N_l}$ for any $l \in [q]$. Here $\mathbf{w}_{li_l} = (w_{li_l j} : j \in [N_l])^\top \in \mathbb{R}^{N_l}$, $w_{li_l j} \geq 0$ and $\|\mathbf{w}_{li_l}\|_1 = 1$. In addition, let $\mathcal{D}^{(l)}$, $\mathcal{D}^{(lm)}$ and $\tilde{\mathcal{D}}^{(l)}$ be defined as in Lemma 29. Furthermore, let $\mathbf{L}_t^{(l)} = (\mathbf{L}_{t, s_1 s_2}^{(l)} : 0 \leq s_1, s_2 \leq t)$ and $\mathcal{L}_{t, s_1 s_2}^{(l)} = (\mathbf{1}_{N_1}^\top \otimes \cdots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \cdots \otimes \mathbf{1}_{N_q}^\top) \mathbf{L}_{t, s_1 s_2}^{(l)} (\mathbf{1}_{N_1} \otimes \cdots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \otimes \cdots \otimes \mathbf{1}_{N_q}) \in \mathbb{R}^{N_l \times N_l}$. Define $\mathcal{L}_t^{(l)} = (\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} : 0 \leq s_1, s_2 \leq T)$ for $l \in [q]$, where $\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} = \mathbf{0}$ for $s_1 > t$ or $s_2 > t$, and $\tilde{\mathcal{L}}_{t, s_1 s_2}^{(l)} = \mathcal{L}_{t, s_1 s_2}^{(l)}$ otherwise. Then we have for any $l \in [q]$,

$$\|\mathcal{D}^{(l)}\|_F^2 \leq c_1 T^{-1}, \quad \|\mathcal{D}^{(l)}\| \leq c_1^{1/2} T^{-1/2}, \quad (\text{A.98})$$

$$\|\mathcal{D}^{(lm)}\|_F^2 \leq c_1 T^{-1}, \quad \|\mathcal{D}^{(lm)}\| \leq c_1^{1/2} T^{-1/2}, \quad (\text{A.99})$$

$$\|\mathbf{D}\|_F^2 \leq c_2 T^{-1}, \quad \|\mathbf{D}\| \leq c_2^{1/2} T^{-1/2}, \quad (\text{A.100})$$

$$\|\tilde{\mathcal{D}}^{(l)}\|_F^2 \leq c_2 T^{-1}, \quad \|\tilde{\mathcal{D}}^{(l)}\| \leq c_2^{1/2} T^{-1/2}, \quad (\text{A.101})$$

$$\max_t \text{tr}(\mathcal{L}_t^{(l)}) < \infty, \quad \text{tr}(\mathbf{L}_t^{(l)}) < \infty, \quad (\text{A.102})$$

where \mathbf{D} is defined in (A.76), and c_1, c_2 are two positive constants.

Proof 1. Proof of (A.98)–(A.99)

Note that by Lemma 31 we have

$$\begin{aligned} \|\mathcal{D}^{(l)}\|_F^2 &= \sum_{s_1, s_2=0}^T \|\mathcal{D}_{s_1 s_2}^{(l)}\|_F^2 \\ &= \sum_{s_1, s_2=0}^T \left\| \left(\mathbf{1}_{N_1}^\top \otimes \cdots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \cdots \otimes \mathbf{1}_{N_q}^\top \right) \mathbf{D}_{s_1 s_2} \right. \\ &\quad \left. \left(\mathbf{1}_{N_1} \otimes \cdots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \otimes \cdots \otimes \mathbf{1}_{N_q} \right) \right\|_F^2 \\ &\leq \sum_{s_1, s_2=0}^T \left(\mathbf{1}_{\prod_l N_l}^\top \mathbf{D}_{s_1 s_2} \mathbf{1}_{\prod_l N_l} \right)^2 \\ &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \left(\sum_{t=\max\{1, s_1, s_2\}}^T \kappa_{\max}^{2t-s_1-s_2} \mathbf{1}_{\prod_l N_l}^\top \mathbf{w}^{(l)} \mathbf{1}_{\prod_l N_l} \right)^2 \\ &= \frac{1}{T^2} \sum_{s_1, s_2=0}^T \left(\sum_{t=\max\{1, s_1, s_2\}}^T \kappa_{\max}^{2t-s_1-s_2} \right)^2 \leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \left(\frac{1}{1-\kappa_{\max}^2} \kappa_{\max}^{|s_1-s_2|} \right)^2 \leq c_1 T^{-1} \end{aligned}$$

where c_1 is a finite constant. Furthermore, noting that $\|\mathcal{D}^{(l)}\| \leq \|\mathcal{D}^{(l)}\|_F \leq c_1^{1/2} T^{-1/2}$, we have the upper bound for $\|\mathcal{D}^{(l)}\|$. Subsequently, the upper bounds for $\|\mathcal{D}^{(lm)}\|_F^2$ and $\|\mathcal{D}^{(lm)}\|$ can be established using the same technique.

2. Proof of (A.100)

First, for any two matrices $\mathbf{M}_1, \mathbf{M}_2$,

$$|\mathbf{M}_1 \mathbf{M}_2^\top|_e \preceq \max_{i,j} |\mathbf{M}_{1i} \mathbf{M}_{2j}^\top| \mathbf{1} \mathbf{1}^\top \leq \max_{i,j} \|\mathbf{M}_{1i}\| \|\mathbf{M}_{2j}\| \mathbf{1} \mathbf{1}^\top.$$

When $\|\mathbf{M}_1\|_{\max} \leq 1$ and $\|\mathbf{M}_2\|_{\max} \leq 1$, we further have $\|\mathbf{M}_{1i}\| \|\mathbf{M}_{2j}\| \leq \|\mathbf{M}_1\|_{\infty} \|\mathbf{M}_2\|_{\infty}$. Hence $|\mathbf{M}_1 \mathbf{M}_2^\top|_e \preceq \|\mathbf{M}_1\|_{\infty} \|\mathbf{M}_2\|_{\infty} \mathbf{1} \mathbf{1}^\top$, and therefore

$$\mathbf{w}^{(l)\top} \mathbf{M}_1 \mathbf{M}_2^\top \mathbf{w}^{(l)} \leq \|\mathbf{M}_1\|_{\infty} \|\mathbf{M}_2\|_{\infty} \mathbf{w}^{(l)\top} \mathbf{1} \mathbf{1}^\top \mathbf{w}^{(l)} = \|\mathbf{M}_1\|_{\infty} \|\mathbf{M}_2\|_{\infty}. \quad (\text{A.103})$$

Note that Lemma 31 implies $\|\mathbf{B}_0^k\|_{\infty} \leq \kappa_{\max}^n < 1$ for any $k > 0$, hence

$$\begin{aligned} \|\mathbf{D}\|_F^2 &= \sum_{s_1, s_2=0}^T \|\mathbf{D}_{s_1 s_2}\|_F^2 \\ &= \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \mathbf{w}^{(l)\top} \mathbf{B}_0^{t_1-s_2} (\mathbf{B}_0^{t_2-s_2})^\top \mathbf{w}^{(l)} \mathbf{w}^{(l)\top} \mathbf{B}_0^{t_2-s_1} (\mathbf{B}_0^{t_1-s_1})^\top \mathbf{w}^{(l)} \\ &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \left(\|\mathbf{B}_0^{t_1-s_2}\|_{\infty} \|\mathbf{B}_0^{t_2-s_2}\|_{\infty} \|\mathbf{B}_0^{t_2-s_1}\|_{\infty} \|\mathbf{B}_0^{t_1-s_1}\|_{\infty} \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \kappa_{\max}^{2t_1+2t_2-2s_1-2s_2} \\
 &\leq \frac{1}{T^2(1-\kappa_{\max}^2)^2} \sum_{s_1, s_2=0}^T \kappa_{\max}^{2|s_1-s_2|} \leq c_2 T^{-1}
 \end{aligned}$$

where c_2 is a finite constant and we used Lemma 31 in the second inequality.

3. Proof of (A.101)

We have

$$\begin{aligned}
 \|\tilde{\mathcal{D}}^{(l)}\|_F^2 &= \sum_{s_1, s_2=0}^T \|\tilde{\mathcal{D}}_{s_1 s_2}^{(l)}\|_F^2 = \sum_{s_1, s_2=0}^T \left\| (\mathbf{1}_{N_1}^\top \otimes \cdots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}}^\top \otimes \cdots \otimes \mathbf{1}_{N_q}^\top) \mathbf{D}_{s_1 s_2} \right\|_F^2 \\
 &= \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \mathbf{w}^{(l)\top} \mathbf{B}_0^{t_1-s_2} \mathbf{B}_0^{t_2-s_2\top} \mathbf{w}^{(l)} \mathbf{w}^{(l)\top} \mathbf{B}_0^{t_2-s_1} \\
 &\quad (\mathbf{1}_{N_1} \mathbf{1}_{N_1}^\top \otimes \cdots \otimes \mathbf{1}_{N_{l-1}} \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \mathbf{1}_{N_{l+1}}^\top \otimes \cdots \otimes \mathbf{1}_{N_q} \mathbf{1}_{N_q}^\top) (\mathbf{B}_0^{t_1-s_1})^\top \mathbf{w}^{(l)} \\
 &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \|\mathbf{B}_0^{t_1-s_2}\|_\infty \|\mathbf{B}_0^{t_2-s_2\top}\|_\infty \\
 &\quad \mathbf{w}^{(l)\top} |\mathbf{B}_0^{t_2-s_1}|_e (\mathbf{1}_{\prod_l N_l} \mathbf{1}_{\prod_l N_l}^\top) |\mathbf{B}_0^{t_1-s_1}|_e^\top \mathbf{w}^{(l)} \\
 &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \|\mathbf{B}_0^{t_1-s_2}\|_\infty \|\mathbf{B}_0^{t_2-s_2}\|_\infty \|\mathbf{B}_0^{t_2-s_1}\|_\infty \|\mathbf{B}_0^{t_1-s_1}\|_\infty \mathbf{w}^{(l)\top} \mathbf{1} \mathbf{1}^\top \mathbf{w}^{(l)} \\
 &= \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \|\mathbf{B}_0^{t_1-s_2}\|_\infty \|\mathbf{B}_0^{t_2-s_2}\|_\infty \|\mathbf{B}_0^{t_2-s_1}\|_\infty \|\mathbf{B}_0^{t_1-s_1}\|_\infty \\
 &\leq \frac{1}{T^2} \sum_{s_1, s_2=0}^T \sum_{t_1, t_2=\max\{1, s_1, s_2\}}^T \kappa_{\max}^{2t_1+2t_2-2s_1-2s_2} \\
 &\leq \frac{1}{T^2(1-\kappa_{\max}^2)^2} \sum_{s_1, s_2=0}^T \kappa_{\max}^{2|s_1-s_2|} \leq c_2 T^{-1}
 \end{aligned}$$

In addition, it further implies that $\|\tilde{\mathcal{D}}^{(l)}\| \leq c_2^{1/2} T^{-1/2}$.

4. Proof of (A.102)

Note that by Lemma 31 we have

$$\begin{aligned}
 \text{tr}(\mathcal{L}_t^{(l)}) &= \sum_{s=0}^t \text{tr}(\mathcal{L}_{t,ss}^{(l)}) \\
 &= \sum_{s=0}^t \text{tr} \left((\mathbf{1}_{N_1}^\top \otimes \cdots \otimes \mathbf{1}_{N_{l-1}}^\top \otimes \mathbf{I}_{N_l} \otimes \cdots \otimes \mathbf{1}_{N_q}^\top) \mathbf{L}_{t,ss}^{(l)} (\mathbf{1}_{N_1} \otimes \cdots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{I}_{N_l} \otimes \cdots \otimes \mathbf{1}_{N_q}) \right) \\
 &\leq \sum_{s=0}^t \left(\mathbf{1}_{\prod_l N_l}^\top |\mathbf{L}_{t,ss}^{(l)}|_e \mathbf{1}_{\prod_l N_l} \right) \leq \sum_{s=0}^t \left(\kappa_{\max}^{2t-2s} \mathbf{1}_{\prod_l N_l}^\top |\mathbf{w}^{(l)}|_e |\mathbf{w}^{(l)\top}|_e \mathbf{1}_{\prod_l N_l} \right) \leq \frac{1}{1-\kappa_{\max}^2} < \infty.
 \end{aligned}$$

Consequently the first inequality of (A.102) holds. Next we note

$$\begin{aligned} \text{tr}(\mathbf{L}_t^{(l)}) &\leq \sum_{s=0}^t (\mathbf{1}_{\prod_l N_l}^\top |\mathbf{L}_{t,ss}^{(l)}|_e \mathbf{1}_{\prod_l N_l}) \\ &\leq \sum_{s=0}^t (\kappa_{\max}^{2t-2s} \mathbf{1}_{\prod_l N_l}^\top |\mathbf{w}^{(l)}|_e |\mathbf{w}^{(l)\top}|_e \mathbf{1}_{\prod_l N_l}) \leq 1/(1 - \kappa_{\max}^2) < \infty, \end{aligned}$$

then the second inequality holds. \blacksquare

Lemma 31. *Under Assumption 5, we have $\|\mathbf{B}_0^n\|_\infty \leq \kappa_{\max}^n$ and $|\mathbf{B}_0|_e^n \mathbf{1}_{\prod_l N_l} \preceq \kappa_{\max}^n \mathbf{1}_{\prod_l N_l}$.*

Proof Note that we have

$$|\mathbf{B}_0|_e \mathbf{1}_{\prod_l N_l} \preceq \sum_l \mathbf{1}_{N_1} \otimes \cdots \otimes \mathbf{1}_{N_{l-1}} \otimes \mathbf{L}_{l,0} \mathbf{1}_{N_l} \otimes \mathbf{1}_{N_{l+1}} \otimes \cdots \otimes \mathbf{1}_{N_q} + |\text{vec}(\mathcal{A}_0)|_e \preceq \kappa_{\max} \mathbf{1}_{\prod_l N_l}$$

by Assumption 5. As a result, we have $\|\mathbf{B}_0^n\|_\infty \leq \| |\mathbf{B}_0|_e^n \|_\infty = \| |\mathbf{B}_0|_e^n \mathbf{1}_{\prod_l N_l} \|_{\max} \leq \kappa_{\max}^n$. \blacksquare

Lemma 32. *Under Assumptions 4 and 5, we have $\|\mathbf{\Gamma}\|_{\max} \leq c_\Gamma$, where $c_\Gamma > 0$ is a constant.*

Proof Define $\mathbf{h}_t = \mathbf{c}_t + \mathbb{E}_t$ and then \mathbf{h}_t follows K' -convex concentration property by Assumption 4 for some constant K' . In addition, let $\mathbb{H}_t = (\mathbf{h}_t^\top, \mathbf{h}_{t-1}^\top, \dots, \mathbf{h}_0^\top)^\top$, $\mathbb{B} = (\mathbf{1}_{\prod_l N_l}, \mathbf{B}_0, \mathbf{B}_0^2, \dots, \mathbf{B}_0^t)$. Note that we have

$$\mathbf{\Gamma} = \text{cov}(\mathbb{Y}_t) = \sum_{k_1, k_2=0}^t \mathbf{B}_0^{k_1} \text{cov} \left\{ \mathbf{h}_{t-k_1}, \mathbf{h}_{t-k_2} \right\} \mathbf{B}_0^{k_2 \top} = \mathbb{B} \text{cov}(\mathbb{H}_t) \mathbb{B}^\top.$$

Then we have $\|\mathbf{\Gamma}\|_{\max} = \max_{i,j} |\mathbf{e}_i^\top \mathbf{\Gamma} \mathbf{e}_j| = \max_i |\mathbf{e}_i^\top \mathbf{\Gamma} \mathbf{e}_i| \leq \|\text{cov}(\mathbb{H}_t)\| \max_i |\mathbf{e}_i^\top \mathbb{B} \mathbb{B}^\top \mathbf{e}_i|$, where $\mathbf{e}_i \in \mathbb{R}^{\prod_l N_l}$ is a vector with its i th element being equal to 1 while others being 0. By Lemma 28, it holds $\|\text{cov}(\mathbb{H}_t)\| \max_i |\mathbf{e}_i^\top \mathbb{B} \mathbb{B}^\top \mathbf{e}_i| \leq 2K^{*2} |\mathbf{e}_i^\top \mathbb{B} \mathbb{B}^\top \mathbf{e}_i|$ for some constant K^* . Next, we have $\mathbf{e}_i^\top \mathbb{B} \mathbb{B}^\top \mathbf{e}_i = \sum_{k=0}^t \mathbf{e}_i^\top \mathbf{B}_0^k \mathbf{B}_0^{k \top} \mathbf{e}_i \leq \sum_{k=0}^t \|\mathbf{B}_0^k\|_\infty^2 \leq \sum_{k=0}^t \kappa_{\max}^{2k} \leq 1/(1 - \kappa_{\max}^2) < \infty$ by Lemma 31 and Assumption 5. Thus the result holds with $c_\Gamma = 1/(1 - \kappa_{\max}^2)$. \blacksquare

Lemma 33. *Under Assumptions 4 and 5, we have $\tau_{\max} = \max_{i_1, \dots, i_q} \lambda_{\max}(\mathbf{\Sigma}_{i_1 \dots i_q}) < \infty$.*

Proof Recall that we have defined

$$\begin{aligned} \mathcal{X}_{i_1 \dots i_q, t} &\stackrel{\text{def}}{=} \left(\sum_{k=1}^{N_1} w_{i_1 k}^{(1)} Y_{ki_2 \dots i_q, (t-1)}, \mathbf{x}_{i_1 t}^{(1)\top}, \dots, \right. \\ &\quad \left. \sum_{k=1}^{N_q} w_{i_q k}^{(q)} Y_{i_1 \dots i_{(q-1)} k, (t-1)}, \mathbf{x}_{i_q t}^{(q)\top}, Y_{i_1 \dots i_q, (t-1)} \right)^\top \in \mathbb{R}^{\sum_l (p_l+1)}. \end{aligned}$$

Note that there are $(2q+1)$ elements in $\mathcal{X}_{i_1 \dots i_q, t}$, denote the m th element as $\mathcal{X}_{i_1 \dots i_q, t}^m$, $1 \leq m \leq (2q+1)$. Specifically, $\mathcal{X}_{i_1 \dots i_q, t}^m = \sum_{k=1}^{N_1} w_{i_l k}^{(l)} Y_{i_1 \dots i_{l-1} k i_{l+1} \dots i_q, (t-1)}$ when m takes $\{1, 3, \dots, 2q-1\}$, $\mathcal{X}_{i_1 \dots i_q, t}^m = \mathbf{x}_{i_l t}^{(l)\top}$ when m takes $\{2, 4, \dots, 2q\}$, and $\mathcal{X}_{i_1 \dots i_q, t}^{2q+1} = Y_{i_1 \dots i_q, (t-1)}$. We write $\boldsymbol{\eta} = (\eta_1, \eta_2, \eta_3, \eta_4, \dots, \eta_{2q+1})^\top$ for any vector $\boldsymbol{\eta} \in \mathbb{R}^{\sum_l (p_l+1)+1}$. Here the even elements, i.e., $\eta_2, \eta_4, \dots, \eta_{2q}$ are vectors with dimension p_2, p_4, \dots, p_{2q} , while other elements are scalar. Then we have

$$\begin{aligned} \lambda_{\max}(\boldsymbol{\Sigma}_{i_1 \dots i_q}) &= \sup_{\|\boldsymbol{\eta}\|=1} \sum_{k_1, k_2=1}^{2q+1} \eta_{k_1}^\top \text{cov}(\mathcal{X}_{i_1 \dots i_q, t}^{k_1}, \mathcal{X}_{i_1 \dots i_q, t}^{k_2}) \eta_{k_2} \\ &\leq (2q+1) \sup_{\|\boldsymbol{\eta}\|=1} \sum_{k=1}^{2q+1} \eta_k^\top \text{var}(\mathcal{X}_{i_1 \dots i_q, t}^k) \eta_k. \end{aligned}$$

Note that $\mathcal{X}_{i_1 \dots i_q, t}^k$ has two typical forms. The first can be written as $\mathbf{w}^\top \mathbb{Y}_t$ ($k = 1, 3, 5, \dots$) by taking \mathbf{w} as a vector with non-negative elements and $\|\mathbf{w}\|_1 = 1$. The second is $\mathbf{x}_{i_l t}^{(l)}$, which represents the covariate information. The first form yields $\eta_k \text{var}(\mathbf{w}^\top \mathbb{Y}_t) \eta_k = \eta_k^2 \mathbf{w}^\top \boldsymbol{\Gamma} \mathbf{w} \leq \|\boldsymbol{\Gamma}\|_{\max} |\mathbf{w}^\top \mathbf{1} \mathbf{1}^\top \mathbf{w}| = \|\boldsymbol{\Gamma}\|_{\max} < c_\Gamma$, where c_Γ is defined in Lemma 32. The second form yields $\eta_k^\top \text{var}(\mathbf{x}_{i_l t}^{(l)}) \eta_k \leq \text{var}(\mathbf{x}_t^{(l)\eta}) < 2K^2$ due to Assumption 4 and Lemma 28. Thus, $\tau_{\max} \leq (2q+1)((q+1)c_\Gamma + qK^2) < \infty$. \blacksquare

Appendix L. Additional Simulation Studies

L.1 Group Number Selection Consistency

To examine the finite sample performance of the group selection consistency, we conduct experiment when $q = 2$. The data generating mechanism is the same as in Section 6.1 in the main text. Specifically, we define

$$\varrho(G_1) = R^{-1} \sum_{r=1}^R I(\widehat{G}_1^{(r)} = G_1), \quad \varrho(G_2) = R^{-1} \sum_{r=1}^R I(\widehat{G}_2^{(r)} = G_2),$$

where $\widehat{G}_1^{(r)}$ and $\widehat{G}_2^{(r)}$ are estimated group numbers in the r th replicate. Hence, $\varrho(G_1)$ and $\varrho(G_2)$ evaluate the proportion of the correctly group numbers for G_1 and G_2 . The results are shown in Table A.7, from which one could see that as the sample size increases, the correct group numbers selection percentage is closer to 1, which can illustrate the group numbers selection consistency in Theorem 3 from the finite sample experiment.

L.2 Random Initialization

In the STEP 2 of Algorithm A.3 in Appendix F, we use the k -means clustering for initialization of the group memberships. Alternatively, we also consider using random initialization in STEP 2, and the estimation results are shown in Table A.8. We take $q = 2$ as an example, and set $N_1 = 100, N_2 = 80$. The time lengths vary in $\{40, 80, 150, 200\}$. All the evaluation metrics are the same as Section 6.1 in the main text. The results demonstrate that

Table A.7: The proportion of selected group numbers G_1 and G_2 in $R = 500$ replicates under different settings.

N_1	N_2	T	G_1	G_2	Scenario 1 (SBM)		Scenario 2 (Power-Law)	
					$\varrho(G_1)$	$\varrho(G_2)$	$\varrho(G_1)$	$\varrho(G_2)$
100	80	20	2	2	0.966	0.966	0.838	0.838
			3	3	0.034	0.034	0.162	0.162
			4	4	0.000	0.000	0.000	0.000
		40	2	2	0.134	0.134	0.000	0.000
			3	3	0.866	0.866	1.000	1.000
			4	4	0.000	0.000	0.000	0.000
200	150	20	2	2	0.756	0.756	0.010	0.010
			3	3	0.244	0.244	0.990	0.990
			4	4	0.000	0.000	0.000	0.000
		40	2	2	0.016	0.016	0.000	0.000
			3	3	0.984	0.984	1.000	1.000
			4	4	0.000	0.000	0.000	0.000

k -means initialization yields superior parameter estimation accuracy compared to random initialization under limited temporal lengths. Notably, as the temporal length increases, the RMSE gap between the two initialization methods narrows significantly, i.e., random initialization based procedure converges to the similar performance as the k -means initialization based procedure. This convergence pattern is consistently observed in mis-classification rates (the last two columns in Table A.8). Critically, k -means initialization maintains robust performance regardless of sample size, while random initialization achieves comparable finite-sample performance when sufficient data is available.

L.3 Robustness Evaluation

To show the robustness of our proposed model, we add several simulations in this section, including model mis-specification and heavy-tail noise settings.

L.3.1 MODEL MIS-SPECIFICATION

In this subsection, we consider two model-specification settings. In both settings, we set $(N_1, N_2) \in \{(100, 80), (200, 150)\}$, and the training time length $T \in \{20, 40, 80, 120\}$.

1. Mis-specified network structures.

First, we consider the case that the network weighting matrices $\{\mathbf{W}^{(l)}\}$ are misspecified when $q = 2$. First, we generate the ‘‘correct’’ networks $\mathbf{A}^{(l)}$ by setting $P(a_{ij,l} = 1) = 1/N_l$ and set $a_{ii} = 0$ by convention. Then, we simulate mis-specified networks by setting $P(a'_{ij,l} = 1) = 1/(10N_l)$ for those $a'_{ij,l} = 0, i \neq j$. Consequently, we observe more links than the correctly specified networks. Denote the mis-specified networks by $\mathbf{A}_{\text{mis}}^{(1)}$ and $\mathbf{A}_{\text{mis}}^{(2)}$. Correspondingly, denote the row-normalized weighting matrices as $\mathbf{W}_{\text{mis}}^{(1)}$ and $\mathbf{W}_{\text{mis}}^{(2)}$. Our data is generated by model (1) with the ‘‘correct’’ weighting matrices $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$, while the estimation is conducted using mis-specified weighting matrices $\mathbf{W}_{\text{mis}}^{(1)}$ and $\mathbf{W}_{\text{mis}}^{(2)}$.

2. Non-linear covariates.

Table A.8: RMSEs of estimated parameters and mis-clustering rates when $G_{1,0} = G_{2,0} = 3$ with 300 replicates. The Initialization methods include random and k -means initialization.

G_1	G_2	N_1	N_2	Initialization	T	$\widehat{\lambda}^{(1)}$	$\widehat{\lambda}^{(2)}$	$\widehat{\zeta}^{(1)}$	$\widehat{\zeta}^{(2)}$	$\widehat{\alpha}$	$\widehat{\eta}_1$	$\widehat{\eta}_2$
3	3	100	80	K -means	40	0.0100	0.0071	0.0140	0.0094	0.0153	0.0313	0.0001
					80	0.0054	0.0046	0.0068	0.0067	0.0092	0.0031	0.0001
					150	0.0040	0.0038	0.0047	0.0060	0.0068	0.0000	0.0001
					200	0.0034	0.0032	0.0041	0.0043	0.0057	0.0000	0.0001
					40	0.0157	0.0135	0.0207	0.0135	0.0356	0.0673	0.0124
					80	0.0089	0.0105	0.0112	0.0093	0.0208	0.0296	0.0044
				Random	150	0.0050	0.0039	0.0067	0.0053	0.0084	0.0117	0.0013
					200	0.0043	0.0040	0.0054	0.0054	0.0083	0.0078	0.0036

Table A.9: ReMSPEs of the mis-specified model and the true model in two scenarios.

N_1	N_2	T	Changing Networks		Nonlinear Covariates	
			ReMSPE _{mis}	ReMSPE _{true}	ReMSPE _{mis}	ReMSPE _{true}
100	80	20	0.7097	0.7082	0.9190	0.8924
		40	0.6845	0.6833	0.8758	0.8497
		80	0.6445	0.6435	0.8589	0.8339
		120	0.6290	0.6282	0.8514	0.8260
200	150	20	0.7309	0.7300	0.9257	0.8985
		40	0.6926	0.6921	0.8792	0.8531
		80	0.6741	0.6736	0.8560	0.8305
		120	0.6415	0.6410	0.8489	0.8236

Second, we consider the case when covariate structure is mis-specified. We generate data using the following model,

$$\begin{aligned}
 Y_{ij,t} = & \lambda_{g_i}^{(1)} \sum_{k=1}^{N_1} \frac{a_{ik}^{(1)}}{n_{1i}} Y_{kj,(t-1)} + \lambda_{g_j}^{(2)} \sum_{k=1}^{N_2} \frac{a_{kj}^{(2)}}{n_{2j}} Y_{ik,(t-1)} \\
 & + \alpha_{g_i^{(1)}g_j^{(2)}} Y_{ij,(t-1)} + f(\mathbf{x}_{it}^{(1)})^\top \boldsymbol{\zeta}_{g_i^{(1)}}^{(1)} + f(\mathbf{x}_{jt}^{(2)})^\top \boldsymbol{\zeta}_{g_j^{(2)}}^{(2)} + \varepsilon_{ij,t}, \quad (\text{A.104})
 \end{aligned}$$

where $f(\mathbf{v}) = (f(v_1), \dots, f(v_{p_i}))$ with $f(v_m) = 0.1v_m^3 + \sin(0.1v_m\pi)$. Then we estimate the model using our proposed algorithm A.1, mis-specifying the true covariates $f(\mathbf{x}_{it}^{(1)})$ and $f(\mathbf{x}_{jt}^{(2)})$ with $\mathbf{x}_{it}^{(1)}$ and $\mathbf{x}_{jt}^{(2)}$.

The data generating scheme for other parts of the model is the same as Section 6.1 in the main text. To evaluate the performance of the two mis-specification settings, we calculate the out-of-sample mean square prediction error (MSPE) for the fitted response $\hat{Y}_{ij,t}$, respectively. Specifically, we set the subsequent $T_{test} = T/2$ samples after the training set as the testing set. Denote the predicted response for the testing set as $\hat{Y}_{ij,t}^{te}$. Correspondingly, calculate the MSPE as

$$\text{MSPE} = (N_1 N_2 T_{test})^{-1} \sum_{t=T_{train}+1}^{T_{test}} \sum_{i,j} (\hat{Y}_{ij,t}^{te} - Y_{ij,t})^2.$$

Then we calculate the relative mean square prediction error as $\text{ReMSPE} = \text{MSPE}/\text{MSPE}_0$, where the baseline $\text{MSPE}_0 = (N_1 N_2 T_{train})^{-1} \sum_{t=1}^{T_{train}} (Y_{ij,t} - \mu_{ij,train})^2$, and $\mu_{ij,train} = T_{train}^{-1} \sum_{t=1}^{T_{train}} Y_{ij,t}$. For comparison, we calculate the ReMSPEs for the true models, denoted as $\text{ReMSPE}_{\text{true}}$. Specifically, we use $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ as the true networks in the first scenario, and use $f(\mathbf{x}_{it}^{(1)})$, $f(\mathbf{x}_{jt}^{(2)})$ as the true covariates in the second scenario. The results are shown in Table A.9. Across both scenarios, ReMSPEs for misspecified models exhibit a pronounced decreasing trend as temporal length increases. Notably, prediction errors under misspecification are close to those of true models, demonstrating the robust performance of the proposed GTNAR framework against model deviations.

L.3.2 HEAVY-TAIL RANDOM NOISE

In this subsection, we generate heavy-tailed data to validate the robustness of our proposed methodology. The experimental design specifies $q = 2$ with $G_{1,0} = G_{2,0} = 3$ underlying groups, with network structures generated from Stochastic Block Model (SBM) networks as in Section 6.1. We conduct $R = 300$ replicates for a reliable evaluation. For each replicate, the idiosyncratic noise $\varepsilon_{i_1 i_2, t}$ follows an independent and identical t -distribution with 5 degrees of freedom (i.e., $\varepsilon_{i_1 i_2, t} \stackrel{\text{i.i.d.}}{\sim} t(5)$). All other simulation parameters align with Section 6.1 of the main text. The results are shown in Table A.10. As both the sample sizes (N_1, N_2) and the time length (T) increase, the RMSEs for all parameters decrease significantly. Our proposed estimators closely approach their oracle counterparts, particularly at larger sample sizes. Notably, the group membership errors $(\hat{\eta}_1, \hat{\eta}_2)$ converge rapidly to zero, reaching near-zero values at $T = 40$ for $N_1 \geq 200$, which confirms the high accuracy of the proposed methodology. This demonstrates the GTNAR method’s robustness in maintaining finite sample performance despite heavy-tailed noise setting.

L.4 Computational Cost

As we comment in Remark 5, the memberships update equation (A.25) can be computed in parallel for each inner-layer subject (i_l in the l th layer). Therefore, we implement the algorithm by a multi-core computational scheme. For each layer, the group memberships estimation is conducted in R version 4.4.3, running on a Apple M3 platform. Parallel processing implements 4 CPU cores using the `doParallelSNOW` package (version 1.0.17). Total available hardware resources included 8 physical cores supporting 8 concurrent threads. We report the average computational cost of 10 replicated experiments in Figure A.7. The results demonstrate the computational efficiency of the proposed methodology under varying network sizes and time lengths. The initialization stage (orange dashed line) exhibits near-constant time complexity with stable execution times (minimal fluctuation as T increases), indicating negligible scaling overhead in initialization. The iterative estimation phase (blue solid line) shows linear scaling with time length T . Notably, even at the maximum network size configuration $(N_1 = 300, N_2 = 250)$, total computational costs remain no more than 20 seconds across all time horizons, which highlights the GTNAR method’s suitability for large-scale applications despite increasing data dimensionality.

L.5 Numerical Convergence Analysis

In this subsection, we take $q = 2$ for example to evaluate the numerical convergence rate of our proposed algorithm. Specifically, denote the loss function in the k th iteration as $Q^{[k]}(\hat{\Theta}^{[k]})$, where $\hat{\Theta}^{[k]}$ is the estimators in the k th iteration. The convergence criterion is set as $|Q^{[k+1]}(\hat{\Theta}^{[k+1]}) - Q^{[k]}(\hat{\Theta}^{[k]})| \leq 10^{-6}$. We first show the steps of iteration required for algorithm convergence and the corresponding total loss defined in equation (8) in each iteration step in Figure A.8. Specifically, we set the network sizes $(N_1, N_2) \in \{(100, 80), (200, 150)\}$ and the time lengths $T \in \{20, 40\}$. We repeat the experiments for 50 replicates in each sample size setting. For a baseline comparison, we calculate the oracle loss by using the oracle estimators, which are obtained by setting the group memberships as true values. Figure A.8 demonstrates the highly efficient convergence of our proposed iterative algorithm

Table A.10: RMSEs of estimated parameters when $G_{1,0} = G_{2,0} = 3$ with noise distribution being $t(5)$. The experiments are conducted with 300 replications. The performances are evaluated for different sample sizes N_1, N_2 and the time length T . Results of the oracle scenario (given true group memberships $\mathcal{G}_1^0, \mathcal{G}_2^0$) are provided. The corresponding CPs are shown in the parenthesis.

N_1	N_2	T	$\hat{\lambda}^{(1)}$	$\hat{\lambda}^{(2)}$	$\hat{\zeta}^{(1)}$	$\hat{\zeta}^{(2)}$	$\hat{\alpha}$	$\hat{\lambda}^{(1)\text{or}}$	$\hat{\lambda}^{(2)\text{or}}$	$\hat{\zeta}^{(1)\text{or}}$	$\hat{\zeta}^{(2)\text{or}}$	$\hat{\alpha}^{\text{or}}$	$\hat{\eta}_1$	$\hat{\eta}_2$	
100	80	20	0.0212 (0.7378)	0.0243 (0.8078)	0.0582 (0.5885)	0.0325 (0.8322)	0.0716 (0.6544)	0.0127 (0.9400)	0.0123 (0.9500)	0.0167 (0.9500)	0.0169 (0.9422)	0.0203 (0.9467)	0.2332	0.0569	
		40	0.0113 (0.8656)	0.0077 (0.9489)	0.0223 (0.8185)	0.0118 (0.9463)	0.0189 (0.8726)	0.0087 (0.9478)	0.0077 (0.9511)	0.0077 (0.9511)	0.0116 (0.9541)	0.0118 (0.9470)	0.0139 (0.9507)	0.0642	0
200	150	20	0.0109 (0.7933)	0.0065 (0.9333)	0.0248 (0.7430)	0.0088 (0.9430)	0.0178 (0.8189)	0.0067 (0.9389)	0.0064 (0.9356)	0.0085 (0.9511)	0.0088 (0.9444)	0.0102 (0.9448)	0.0829	0.0001	
		40	0.0054 (0.9133)	0.0043 (0.9356)	0.0074 (0.9256)	0.0060 (0.9563)	0.0075 (0.9381)	0.0049 (0.9367)	0.0043 (0.9333)	0.0060 (0.9485)	0.0060 (0.9485)	0.0060 (0.9563)	0.0069 (0.9519)	0.0068	0.0001
300	250	20	0.0053 (0.9133)	0.0039 (0.9356)	0.0085 (0.9256)	0.0053 (0.9563)	0.0077 (0.9381)	0.0042 (0.9367)	0.0039 (0.9333)	0.0055 (0.9485)	0.0055 (0.9485)	0.0053 (0.9563)	0.0064 (0.9519)	0.0153	0
		40	0.0027 (0.9600)	0.0027 (0.9467)	0.0039 (0.9478)	0.0038 (0.9437)	0.0044 (0.9463)	0.0027 (0.9600)	0.0027 (0.9467)	0.0027 (0.9478)	0.0039 (0.9478)	0.0038 (0.9437)	0.0044 (0.9467)	0	0

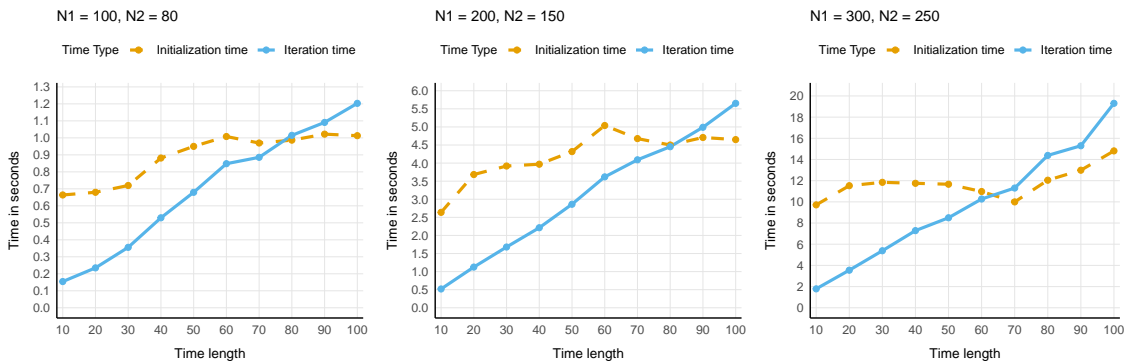


Figure A.7: Computational costs (in second) of the initialization (orange dashed line) and iterative estimation (blue solid line) stage under different sample sizes and time lengths.

to the oracle loss across diverse sample size configurations. Notably, either larger network sizes ($N_1 = 200, N_2 = 150$) or longer time periods ($T = 40$) enhance convergence speed.

We further show the average and the maximum number of iterations. Fix network sizes as $N_1 = 100, N_2 = 80$, we show the results as the time length T grows from 10 to 100; and fix the time length $T = 20$, we increase the network sizes throughout $(N_1, N_2) \in \{(30, 20), (50, 40), (80, 60), (100, 80), (120, 100), (160, 120), (200, 150), (250, 180), (300, 250)\}$. The results are shown in Figure A.9. The average iteration number demonstrates consistent monotonic decline as sample sizes increase, which highlights the algorithm’s intrinsic efficiency. While the maximum iteration number exhibits fluctuations at minimal network sizes, it undergoes rapid stabilization followed by progressive decline as sample sizes grow.

L.6 Simulation Results when $q = 3$

In this subsection, we conduct experiment when $q = 3$ to show the finite sample performance of our proposed method. We set the true group numbers in the three dimensions as $G_{1,0} = 3, G_{2,0} = 3, G_{3,0} = 2$. The true parameters are provided in Table A.11. We consider the networks generated from stochastic block models. The network sizes are set as $(N_1, N_2, N_3) \in \{(20, 20, 20), (30, 30, 30)\}$, and the time length is set to be $T \in \{10, 40, 80\}$. The network and the data generation schemes are the same as those in Section 6.1 in the main text. The experiments are repeated for $R = 300$ times.

We first evaluate estimation accuracy under correctly specified group numbers, with comprehensive results presented in Table A.12. We observe three key patterns. First, the mis-classification rates decline rapidly to zero as sample sizes increase, confirming precise recovery of group memberships. Second, the RMSEs decrease monotonically with increasing network sizes or time length, demonstrating consistent estimation efficiency. Third, the coverage probabilities converge to the nominal 95% confidence level under larger samples, validating asymptotic normality of the estimators. Collectively, these findings meet theoretical results for statistical consistency.

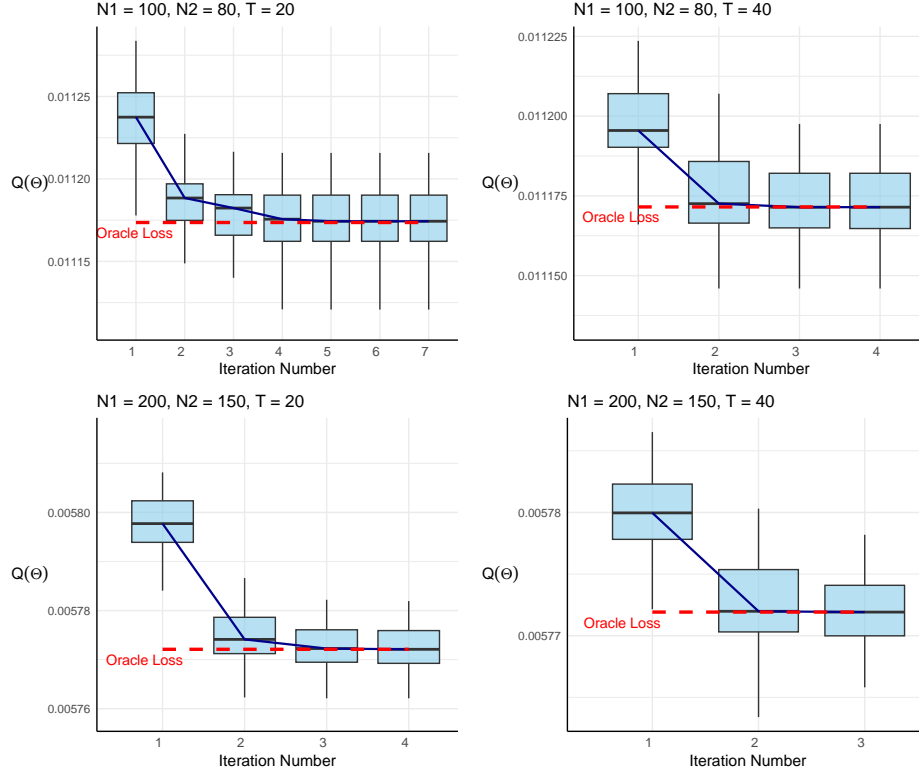


Figure A.8: Total loss in each iteration step and the number of iterations required for algorithm convergence in different network sizes and time lengths. Each box shows the total loss of $R = 50$ replicates. The red line shows the oracle loss.

Table A.11: True parameters when $G_{1,0} = 3, G_{2,0} = 3, G_{3,0} = 2$.

$G_{1,0} = 3$	$G_{2,0} = 3$	$G_{3,0} = 2$
$\lambda_{g^{(1)}}^{(1)}$ (-0.1, 0.2, 0.3)	$\lambda_{g^{(2)}}^{(2)}$ (0.15, 0.2, 0.4)	$\lambda_{g^{(3)}}^{(3)}$ (-0.2, 0.25)
$\zeta_{g^{(1)}}^{(1)}$ $\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \\ 0.24 & 0.30 & -0.32 \end{pmatrix}$	$\delta_{g^{(2)}}^{(2)}$ $\begin{pmatrix} 0.25 & -0.3 & 0.35 \\ 0.2 & -0.25 & 0.32 \\ 0.1 & -0.2 & 0.2 \end{pmatrix}$	$\delta_{g^{(3)}}^{(3)}$ $\begin{pmatrix} -0.1 & 0.2 & 0.4 \\ 0.3 & 0.1 & -0.32 \end{pmatrix}$
$\alpha_{..1}$ $\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \\ 0.24 & 0.30 & -0.32 \end{pmatrix}$	$\alpha_{..2}$ $\begin{pmatrix} 0.2 & 0.25 & -0.3 \\ 0.15 & 0.35 & -0.35 \\ 0.24 & 0.30 & -0.32 \end{pmatrix}$	

Table A.12: RMSEs of estimated parameters when $G_{1,0} = G_{2,0} = 3$ and $G_{3,0} = 2$ with 300 replications. The performances are evaluated for different sample sizes N_1, N_2, N_3 and the time length T . The corresponding CPs are shown in the parenthesis.

N_1	N_2	N_3	T	$\hat{\lambda}^{(1)}$	$\hat{\lambda}^{(2)}$	$\hat{\lambda}^{(3)}$	$\hat{\zeta}^{(1)}$	$\hat{\zeta}^{(2)}$	$\hat{\zeta}^{(3)}$	$\hat{\alpha}$	$\hat{\eta}_1$	$\hat{\eta}_2$	$\hat{\eta}_3$
20	20	20	10	0.0165 (0.9256)	0.0141 (0.9411)	0.0097 (0.9367)	0.0207 (0.9226)	0.0204 (0.9315)	0.0128 (0.9272)	0.0666 (0.9222)	0.0083	0	0
			40	0.0066 (0.9356)	0.0067 (0.9378)	0.0044 (0.9533)	0.0093 (0.9411)	0.0096 (0.9452)	0.0059 (0.9372)	0.0263 (0.9422)	0	0	0
30	30	30	10	0.0025 (0.9311)	0.0024 (0.9467)	0.0024 (0.9383)	0.0065 (0.9459)	0.0065 (0.9519)	0.0042 (0.9478)	0.0158 (0.9478)	0	0	0
			40	0.0105 (0.9422)	0.0082 (0.9400)	0.0057 (0.9383)	0.0130 (0.9196)	0.0106 (0.9370)	0.0068 (0.9378)	0.0435 (0.9269)	0.0078	0	0
80	80	80	10	0.0034 (0.9489)	0.0033 (0.9567)	0.0023 (0.9583)	0.0050 (0.9526)	0.0049 (0.9556)	0.0033 (0.9411)	0.0133 (0.9459)	0	0	0
			40	0.0013 (0.9478)	0.0013 (0.9422)	0.0012 (0.9650)	0.0034 (0.9489)	0.0037 (0.9507)	0.0024 (0.9461)	0.0084 (0.9463)	0	0	0

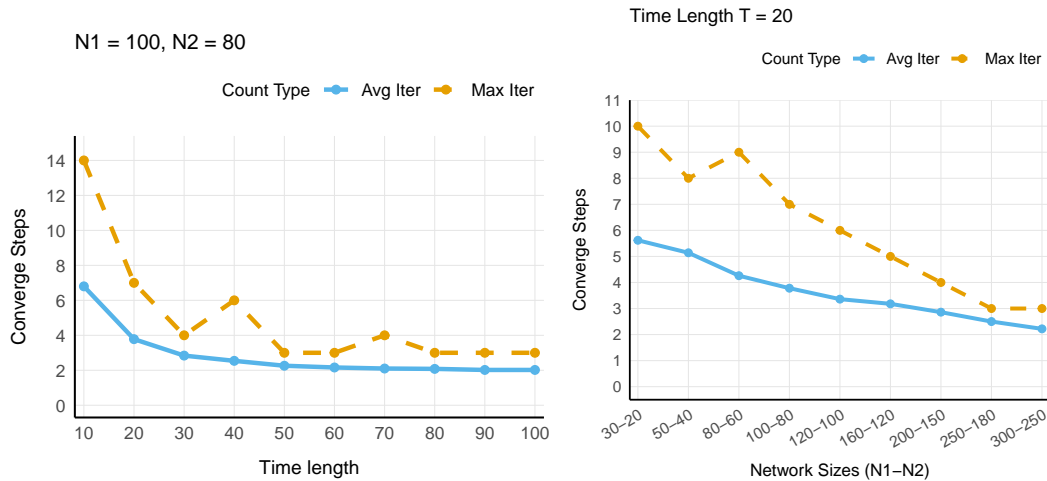


Figure A.9: Average and maximum number of iterations required for convergence under different sample sizes throughout 50 replicates. The orange dashed line shows the maximum number, while the blue solid lines shows the average one.

Next, we consider the setting of unspecified group numbers. We implement the QIC introduced in Section 3.3 in the main text, with κ set as the same with Section 6.1. The results are shown in Table A.13. When the group numbers are under-specified (i.e., $G_l = 2$, $l = 1, 2, 3$), the node-wise RMSEs are large and do not decrease as the sample sizes increase. On the contrary, when the group numbers are correctly- or over-specified, the node-wise RMSEs show a decreasing pattern as the theoretical results show. These findings are consistent with those for two dimensional tensor in Section 6.2 of the main text.

L.7 Comparison with Existing Methods

L.7.1 COMPARING WITH ESTIMATION USING A SERIES OF MATRICES MODELS

To enhance the difference between our proposed GTNAR model and a series of matrix models, we conduct an experiment when $q = 3$ to show the fundamental role of the intrinsic inner-tensor data information. Specifically, we generate data using the same setting as in Appendix L.6. The total time length is set to be $T_{train} + T_{test}$, and we use the first T_{train} as training set, while leaving the following T_{test} as the testing set. By implementing Algorithm A.2, we obtain the in-sample and out-of-sample predicted values $\hat{\mathbf{y}}_t^{\text{GTNAR}}$. For comparison, we divide the data by the third dimension into N_3 parallel data slices. Each slice contains a matrix-valued time series. By algorithm A.1, we can estimate the two-dimensional parameters for each series slice. Take the i_3 slice for example, we can calculate the in-sample fitted matrix $\hat{\mathbf{Y}}_{\cdot i_3, t}$, ($t \in [T_{train}]$) and out-of-sample prediction values $\hat{\mathbf{Y}}_{\cdot i_3, t}$, ($t = T_{train} + 1, \dots, T_{train} + T_{test}$) for each slice. Then we aggregate the fitted and predicted values in all N_3 slices, obtaining the in-sample and out-of-sample predicted tensor $\hat{\mathbf{y}}_t^{\text{MatSlice}}$. We call this method by ‘‘Matrix Slices’’ method. Using the same metrics as in Appendix L.3.1, we calculate ReMSPE for both methods regarding in-sample and out-of-

Table A.13: Simulation results with pre-specified group numbers as well as the QIC selection group numbers \hat{G}_l , $l = 1, 2, 3$. The true group numbers are set as $G_{1,0} = G_{2,0} = 3$ and $G_{3,0} = 2$. The node-wise RMSEs of different estimators are denoted as $\hat{\lambda}_{\text{all}}^{(l)}$, $\hat{\zeta}_{\text{all}}^{(l)}$, α_{all} for $l = 1, 2, 3$.

N_1	N_2	N_3	T	G_1	G_2	G_3	$\hat{\lambda}_{\text{all}}^{(1)}$	$\hat{\lambda}_{\text{all}}^{(2)}$	$\hat{\lambda}_{\text{all}}^{(3)}$	$\hat{\zeta}_{\text{all}}^{(1)}$	$\hat{\zeta}_{\text{all}}^{(2)}$	$\hat{\zeta}_{\text{all}}^{(3)}$	$\hat{\alpha}_{\text{all}}$	$\hat{\xi}_1$	$\hat{\xi}_2$	$\hat{\xi}_3$			
20	20	20	40	Oracle	2	2	0.0037	0.0038	0.0031	0.0052	0.0046	0.0040	0.0058	-	-	-			
					2	2	0.0559	0.0805	0.0551	0.0361	0.0468	0.0051	0.1400	0.2500	0.2080	0			
					3	3	0.0037	0.0038	0.0031	0.0052	0.0046	0.0040	0.0174	0.1403	0.0142	0			
				\hat{G}_1	\hat{G}_2	\hat{G}_3	4	4	4	0.0047	0.0048	0.0047	0.0061	0.0060	0.0063	0.0123	0	0	0
							4	4	4	0.0037	0.0038	0.0031	0.0052	0.0046	0.0040	0.0058	-	-	-
							Oracle	2	2	0.0014	0.0017	0.0017	0.0034	0.0034	0.0026	0.0037	-	-	-
20	20	20	80	Oracle	2	2	0.1140	0.0571	0.0580	0.0343	0.0213	0.0041	0.1447	0.2449	0.1959	0			
					3	3	0.0037	0.0028	0.0028	0.0040	0.0038	0.0026	0.0065	0	0	0			
					4	4	4	0.0040	0.0033	0.0038	0.0046	0.0042	0.0043	0.0099	0	0	0		
				\hat{G}_1	\hat{G}_2	\hat{G}_3	4	4	4	0.0014	0.0017	0.0017	0.0034	0.0034	0.0026	0.0037	-	-	-
							4	4	4	0.0014	0.0017	0.0017	0.0034	0.0034	0.0026	0.0037	-	-	-
							Oracle	2	2	0.0014	0.0017	0.0017	0.0034	0.0034	0.0026	0.0037	-	-	-

Table A.14: In-sample and out-of-sample ReMSPE of GTNAR and the matrix slices estimation methods.

N_1	N_2	N_3	T	GTNAR		Matrix Slices	
				ReMSPE _{tr}	ReMSPE _{te}	ReMSPE _{tr}	ReMSPE _{te}
20	20	20	10	0.6739	0.6549	0.8836	1.0319
			40	0.5882	0.5871	0.7692	0.8003
			80	0.0234	0.0049	2.4769	4.3342
30	30	30	10	0.6261	0.6173	0.8789	1.0281
			40	0.4710	0.4660	0.9787	1.7442
			80	0.0066	0.0006	4.8092	2.9876

sample performance. The experiments are both repeated for $R = 100$ times. The results are shown in Table A.14. GTNAR exhibits consistent improvement in both in-sample and out-of-sample prediction accuracy as sample size increases. In contrast, Matrix Slices shows no systematic improvement trend, with error magnitudes actually increasing substantially at larger T . Compared with the Matrix Slices method, our proposed GTNAR shows a superior prediction performance. These results show that the GTNAR fundamentally differs from simple matrix slicing approaches. The error gaps reveals critical information loss inherent in slice-based approximations when the inherent data is generated from the tensor-valued model.

L.7.2 COMPARING WITH OTHER METHODS

In this section, we conduct the prediction and compare the accuracy with a number of competing methods. Denote T_{train} and T_{test} as the time length for training set and the testing set. We generate the simulation data when $q = 2$, and the generation scheme is the same as that in Section 6.1 in the main text. We set the sample sizes as $(N_1, N_2, T_{train}) \in \{(30, 20, 20), (50, 40, 30), (80, 60, 40), (100, 80, 50), (120, 100, 60), (160, 120, 70), (200, 150, 80), (250, 180, 90)\}$. For each sample size setting, we set the following $T_{test} = 0.5T_{train}$ as the test set. To evaluate the prediction accuracy, we repeat the experiments for each setting for $R = 50$ times, and we compute average relative mean square prediction error (ReMSPE). Specifically, first calculate mean square prediction error as

$$\text{MSPE}^r = (N_1 N_2 T_{test})^{-1} \sum_{t=t_{test,0}}^{t_{test,0}+T_{test}} \sum_{i,j} (\hat{Y}_{ij,t}^r - Y_{ij,t}^r)^2,$$

where $\hat{Y}_{ij,t}^r$ is the predicted response in the r th replicate, and $t_{test,0} = T_{train} + 1$ is the starting time point. Then, calculate

$$\text{MSPE}_0^r = (N_1 N_2 T_{test})^{-1} \sum_{t=t_{test,0}}^{t_{test,0}+T_{test}} \sum_{i,j} (Y_{ij,t}^r - \hat{\mu}_{ij,train}^r)^2,$$

where $\hat{\mu}_{ij,train}^r = T_{train}^{-1} \sum_{t=1}^{T_{train}} Y_{ij,t}^r$ as the mean response in the training set. Subsequently, we calculate the average ReMSPE as $\text{ReMSPE} = R^{-1} \sum_r \text{MSPE}^r / \text{MSPE}_0^r$.

Specifically, we compare the prediction accuracy estimated by the proposed algorithm with that estimated by the sparse VAR method (sVAR, Nicholson et al. (2020)), the multilinear tensor regression model (BiTR, Hoff (2015)), the group NAR model (GNAR, Zhu et al. (2023)), and the deep learning method for time series, LSTM (Hochreiter and Schmidhuber, 1997). For the sVAR model, we first vectorize the matrix response \mathbf{Y}_t to be in the vector form as $\text{vec}(\mathbf{Y}_t)$. Subsequently, we apply the sVAR method in the `bigtime` package to $\{\text{vec}(\mathbf{Y}_t)\}$ to obtain the model estimation and prediction result. For BiTR model, we use \mathbf{Y}_t as the response variable, and set the explanatory variables as the concatenation by the row network term $\mathbf{W}^{(1)}\mathbf{Y}_{t-1}$, the and column network term $\mathbf{Y}_{t-1}\mathbf{W}^{(2)}$, the lag term \mathbf{Y}_{t-1} and the covariates $\mathbf{X}_t^{(1)}, \mathbf{X}_t^{(2)}$. We apply the alternating least squares algorithm proposed by Hoff (2015) to estimate the model. For the GNAR model proposed by Zhu et al. (2023), we note that it only involves a single network and can only be applied for the vector time series data. Therefore we make prediction with their model form by involving one network matrix a time. To involve the user network (i.e., row network), we apply the GNAR model in the following form

$$Y_{ij,t} = \sum_{k=1}^{N_1} \beta_{g_i^{(1)} g_k^{(1)}} w_{1,ik} Y_{kj,t-1} + \nu_{g_i^{(1)}} Y_{ij,t-1} + \mathbf{x}_i^{(1)\top} \boldsymbol{\zeta}_{g_i^{(1)}} + \varepsilon_{ij,t}, \quad (\text{A.105})$$

and estimate the model parameters and make prediction for each $j \in [N_2]$. The model prediction performance is denoted as GNAR-R in Table A.18. Similarly, to incorporate the spatial network (i.e., column network), we estimate the following GNAR model as

$$Y_{ij,t} = \sum_{k=1}^{N_2} \beta_{g_k^{(2)} g_j^{(2)}} w_{2,kj} Y_{ik,t-1} + \nu_{g_j^{(2)}} Y_{ij,t-1} + \mathbf{x}_j^{(2)\top} \boldsymbol{\zeta}_{g_j^{(2)}} + \varepsilon_{ij,t},$$

for each $i \in [N_1]$, whose prediction performance is denoted as GNAR-C in Figure A.10. For the above two models, we use the open source code provided by Zhu et al. (2023) to obtain the result. Lastly, we apply the LSTM method (Hochreiter and Schmidhuber, 1997) to our data $\{(\mathbf{x}_{it}^{(1)} \in \mathbb{R}^{p_1}, \mathbf{x}_{jt}^{(2)} \in \mathbb{R}^{p_2}, Y_{ij,t}) : i \in [N_1], j \in [N_2], t \in [T]\}$. During training, we set the learning rate as 0.05, the dimension of hidden layers as 64, and set the number of iteration as 100. The momentum coefficient is set as 0.5.

We remark that due to the suboptimal estimation accuracy of BiTR, which exhibits orders of magnitude differences compared to others, we present the ReMSPEs of GTNAR, GNAR-R, GNAR-C, sVAR and LSTM in Figure A.10 for visual clarity. The results for BiTR are provided in Table A.15 for detailed examination. Based on the simulation results, when the data is generated using the GTNAR model, the prediction accuracy of all alternative methods proves suboptimal. Even the LSTM model, which is recognized for its strong fitting capabilities by its non-linear structure, fails to achieve satisfactory prediction accuracy. Furthermore, in Appendix M.3, we conduct rolling-window prediction comparisons across six methods using the Yelp dataset, which also demonstrate the superiority of our proposed approach.

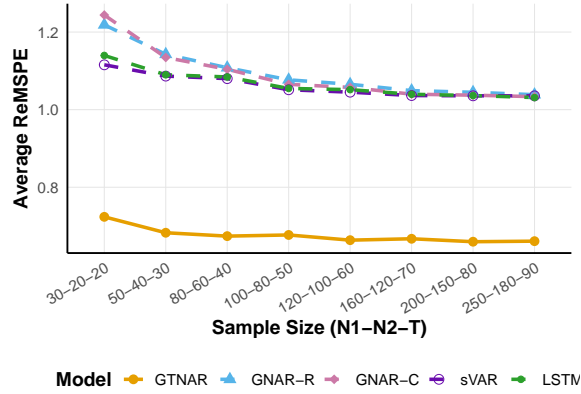


Figure A.10: Average ReMSPEs of GTNAR, GNAR-R, GNAR-C, sVAR and LSTM under different sample sizes.

Table A.15: Average ReMSPEs of GTNAR and BiTR under different sample sizes.

N_1	N_2	T_{train}	GTNAR	BiTR
30	20	20	0.7237	94.8849
50	30	30	0.6827	160.6119
80	60	40	0.6740	219.2674
100	80	50	0.6770	138.1804
120	100	60	0.6636	105.4973
160	120	70	0.6672	139.3628
200	150	80	0.6596	101.1558
250	180	90	0.6610	28.5175

Appendix M. Additional Yelp Data Analysis

M.1 Additional Results for Yelp Data Analysis

We visualize the relationship between these user-related covariates and the response variable in Figure A.11. The plot reveals that users who receive more tags for their reviews tend to be motivated to contribute more reviews. Notably, VIP users in Scottsdale and Toronto tend to write more reviews, whereas VIP users in Charlotte exhibit comparatively less activity. The estimation results for city Scottsdale and Toronto are provided in Table A.16.

M.2 Ablation Experiment

In this subsection, we use ablation experiments to show the necessity of each component in our GTNAR model. We implement both the GTNAR model (1) and the Mixed GTNAR model (25) to analyze the Yelp dataset. The evaluation metric is the similar as that in Appendix L.7.2 yet under a rolling window. Specifically, first calculate mean square prediction error as

$$\text{MSPE} = (N_1 N_2 T_{test})^{-1} \sum_{t=t_{test,0}}^{t_{test,0}+T_{test}} \sum_{i,j} (\hat{Y}_{ij,t} - Y_{ij,t})^2,$$

Then, calculate

$$\text{MSPE}_0 = (N_1 N_2 T_{test})^{-1} \sum_{t=t_{test,0}}^{t_{test,0}+T_{test}} \sum_{i,j} (Y_{ij,t} - \hat{\mu}_{ij,train})^2,$$

where $\hat{\mu}_{ij,train} = T_{train}^{-1} \sum_{t=t_{test,0}-T_{train}}^{t_{test,0}-1} Y_{ij,t}$ as the mean response in the training set. Subsequently, we calculate the ReMSPE as $\text{ReMSPE} = \text{MSPE}/\text{MSPE}_0$. Recall the GTNAR model

$$\mathbf{Y}_t = \underbrace{(\mathbf{L}^{(1)} \mathbf{W}^{(1)}) \mathbf{Y}_{t-1} + \mathbf{Y}_{t-1} (\mathbf{W}^{(2)} \mathbf{L}^{(2)})}_{\text{network terms}} + \underbrace{\mathbf{A} \odot \mathbf{Y}_{t-1}}_{\text{momentum term}} + \underbrace{\boldsymbol{\beta}_{X_1,t}^{(1)} \mathbf{1}_{N_2}^\top + \mathbf{1}_{N_1} \boldsymbol{\beta}_{X_2,t}^{(2)\top}}_{\text{covariates terms}} + \mathbf{E}_t, \quad (\text{A.106})$$

For the ablation experiment, we remove the network terms, the lag term, and the covariates terms from GTNAR (A.106) and perform predictions, respectively. For performance evaluation, we predict the out-of-sample review number in the same rolling window settings, and calculate the ReMSPE as defined in Appendix L.7.2. The prediction results are shown in Table A.17, from which we can see that both GTNAR and Mixed GTNAR consistently achieve the highest prediction accuracy across most of the five evaluation testing time periods. This performance superiority is particularly evident in Charlotte, Las Vegas, and Scottsdale, where these models outperform all ablated versions at every starting time point. The ablation study reveals performance deterioration when each component is removed, validating the contribution of each model component (network terms, momentum term, and covariates terms) to GTNAR’s robust predictive capability.

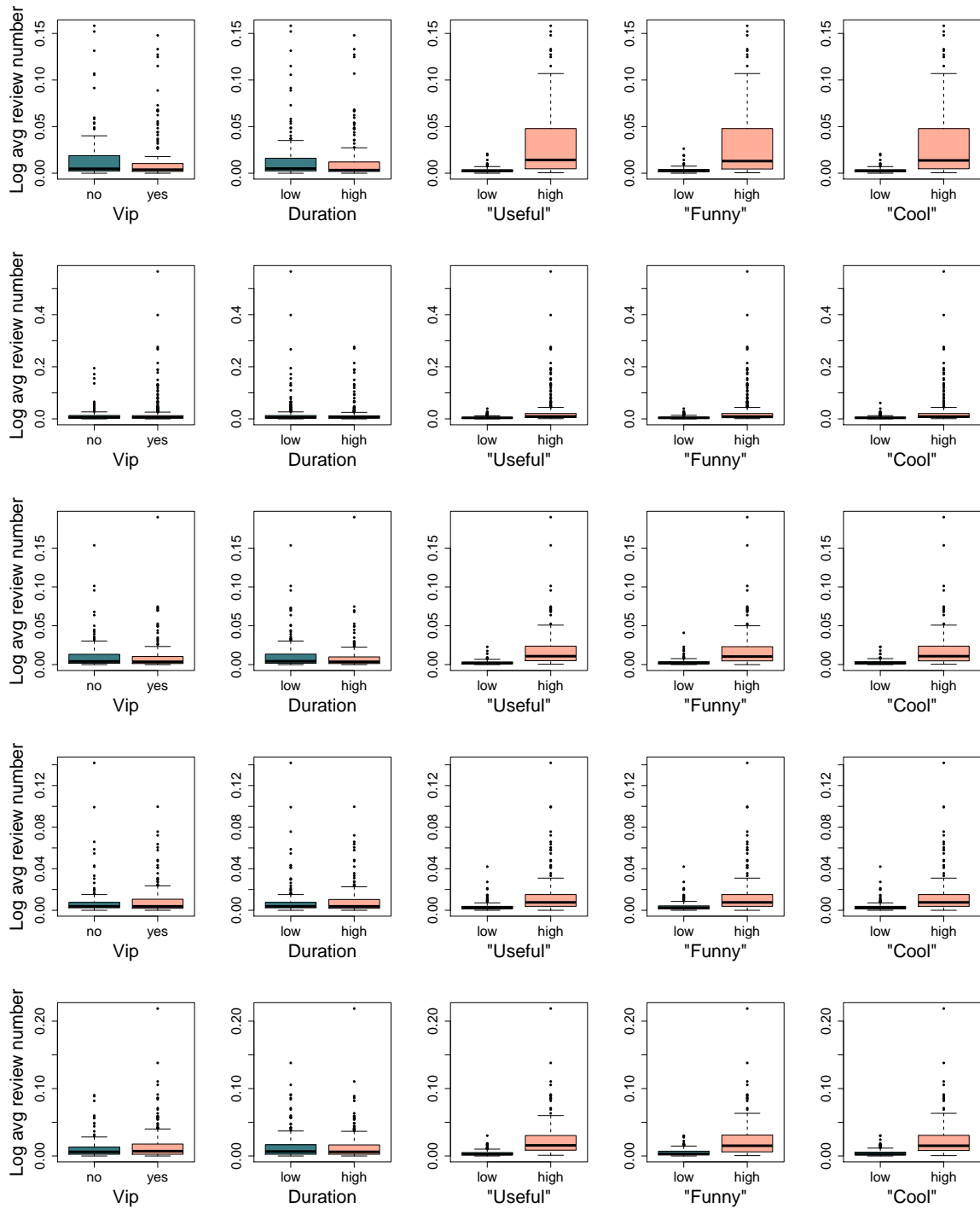


Figure A.11: Boxplots for response variable with regards to the users' covariates in Charlotte, Las Vegas, Phoenix, Scottsdale, and Toronto (displayed in each line). Specifically, the y-axis is the log-transformed average review number given by each user throughout the time span. The x-axis shows the corresponding high or low level for each covariate, separated by the median for continuous variables.

Table A.16: Estimation results for Scottsdale and Toronto. The p -values are shown in the parenthesis.

Parameters	Scottsdale ($N_1 = 391, N_2 = 60$)		Toronto ($N_1 = 462, N_2 = 56$)	
	$\lambda_{g^{(1)}}^{(1)}$	$\lambda_{g^{(2)}}^{(2)}$	$\lambda_{g^{(1)}}^{(1)}$	$\lambda_{g^{(2)}}^{(2)}$
	0.067 (<0.001)	0.011 (<0.001)	0.069 (<0.001)	0.239 (<0.001)
	0.013 (<0.001)	0.037 (<0.001)	0.016 (<0.001)	0.107 (<0.001)
Intercept	$\zeta_{g^{(1)}}^{(1)}$ -0.002 (0.175)	$\zeta_{g^{(2)}}^{(2)}$ 0.002 (<0.001)	$\zeta_{g^{(1)}}^{(1)}$ -0.009 (0.001)	$\zeta_{g^{(2)}}^{(2)}$ 0.016 (<0.001)
$\zeta_{dur}^{(1)} / \zeta_{star}^{(2)}$	0.009 (<0.001)	0.002 (<0.001)	10^{-4} (<0.001)	0.003 (0.274)
$\zeta_{vip}^{(1)} / \zeta_{num}^{(2)}$	0.002 (0.037)	0.001 (<0.001)	0.001 (0.453)	0.003 (0.209)
$\zeta_{use}^{(1)}$	0.227 (<0.001)	0.023 (<0.001)	0.001 (<0.001)	0.003 (0.005)
$\zeta_{fun}^{(1)}$	-0.360 (<0.001)	-0.025 (<0.001)	0.069 (0.155)	0.005 (<0.001)
$\zeta_{cool}^{(1)}$	0.176 (<0.001)	0.028 (<0.001)	0.032 (<0.001)	0.003 (0.009)
	$\alpha^T \in \mathbb{R}^{G_1 \times G_2}$		$\alpha^T \in \mathbb{R}^{G_1 \times G_2}$	
	0.071 (<0.001)	0.023 (<0.001)	0.055 (<0.001)	0.290 (<0.001)
	0.235 (<0.001)	0.051 (<0.001)	0.033 (<0.001)	0.096 (<0.001)

Table A.17: ReMSPEs of the ablation study in five cities under the rolling window setting.

		Starting time point				
		2010-Q1	2010-Q2	2010-Q3	2010-Q4	2011-Q1
Charlotte	GTNAR	0.8386	0.8132	0.8361	0.8646	0.8951
	Mixed GTNAR	0.8228	0.7725	0.8382	0.9019	0.9629
	no network terms	0.8421	0.8140	0.8392	0.8689	0.8979
	no momentum term	0.8717	0.8199	0.8629	0.8824	0.9624
	no covariates terms	0.8539	0.8297	0.8565	0.8772	0.9131
Las Vegas	GTNAR	0.8336	0.7811	0.8549	0.9102	0.9601
	Mixed GTNAR	0.8228	0.7725	0.8382	0.9019	0.9629
	no network terms	0.8373	0.8021	0.8519	0.9237	0.9784
	no momentum term	0.8238	0.7913	0.8520	0.9145	0.9729
	no covariates terms	0.8119	0.7918	0.8395	0.9070	0.9691
Phoenix	GTNAR	1.1625	1.0784	0.9419	0.8730	0.8581
	Mixed GTNAR	1.1406	1.0888	0.9571	0.8822	0.8582
	no network terms	1.1048	1.0793	0.9393	0.8677	0.8616
	no momentum term	1.1426	1.0859	0.9424	0.8731	0.8554
	no covariates terms	1.1310	1.1023	0.9542	0.8877	0.8664
Scottsdale	GTNAR	0.8725	0.8588	0.8535	0.8757	0.8415
	Mixed GTNAR	0.8986	0.8757	0.8629	0.8859	0.8581
	no network terms	0.8761	0.8629	0.8587	0.8771	0.8441
	no momentum term	0.8992	0.8699	0.8545	0.9057	0.8931
	no covariates terms	0.8831	0.8685	0.8636	0.8854	0.8492
Toronto	GTNAR	1.5779	1.4670	1.3785	1.2241	1.0669
	Mixed GTNAR	1.5829	1.5279	1.3961	1.2313	1.0773
	no network terms	1.5751	1.4678	1.3863	1.2334	1.0681
	no momentum term	1.5783	1.4787	1.4079	1.2362	1.0766
	no covariates terms	1.5995	1.4942	1.4019	1.2513	1.0877

M.3 Prediction Accuracy of Comparing Methods

In this section, we conduct the prediction for the Yelp dataset and compare the accuracy with a number of competing methods. To illustrate the prediction performance of our proposed GTNAR model, we take a rolling window prediction setting. Specifically, we set the training length as $T_{train} = 25$ from the first data point, and set $T_{test} = 5$. Then we use a sliding window approach to evaluate the prediction performance by calculating the Root Mean Square Prediction Error (RMSPE) with moving one quarter in each step, which is calculated by

$$\text{RMSPE} = \left\{ (N_1 N_2 T_{test})^{-1} \sum_{t=t_{test,0}}^{t_{test,0}+T_{test}} \sum_{i,j} (\hat{Y}_{ij,t} - Y_{ij,t})^2 \right\}^{1/2},$$

where $\hat{Y}_{ij,t}$ is the predicted response from user i to the district j in the t th quarter, and $t_{test,0}$ is the starting time point. The settings for the competing methods, namely, sVAR, BiTR, GNAR-R, GNAR-C and LSTM are the same as those in Appendix L.7.2.

The prediction results are detailed in Table A.18. In comparison to the sVAR, BiTR, GNAR-R and GNAR-C, the GTNAR method demonstrates superior and more stable prediction accuracy across various starting points in all five cities. In comparison to the LSTM model, we find that the prediction performances are comparable. In particular, for Charlotte, Las Vegas, Scottsdale and Toronto, the GTNAR method outperforms the LSTM with higher prediction accuracy. For Phoenix, the prediction error of the LSTM method in the first start time point is slightly smaller than that of the GTNAR method. We would like to remark that other than our relative robust prediction performance, our modeling approach provides a clearer model interpretation compared with complex nonlinear machine learning methods.

Appendix N. Additional Trading Data Application

N.1 Data Description

In this section, we embark on an analysis of the monthly import and export volumes of goods among countries/regions, spanning the period from 1993 to 2022, comprising a total of $T = 360$ months. The source for our multilateral trading data is the IMF-DOTS database (IMF, 2017). To enhance the stationarity of the time series, we define the response variable, denoted by $Y_{ij,t}$, as the first-order difference in export volume, which is referred to as the “increased export volume (IEV)” in the subsequent discussions. IEV quantifies the change in export volume from the i th country/region to the j th country/region between the t th month and the preceding $(t - 1)$ th month.

Covariates. In our analysis, we consider three essential covariates for each country or region. Firstly, the annual gross domestic product ($x_{it,\text{gdp}}^{(1)}/x_{jt,\text{gdp}}^{(2)}$) serves as a measure of the overall size of the economy. Secondly, we incorporate the difference of log-annual average exchange rate to US dollars ($x_{it,\text{exc}}^{(1)}/x_{jt,\text{exc}}^{(2)}$) in consecutive years, which influences the actual costs of exporting goods. The influence from exchange rate toward trading pattern is underscored in prior research (Hooper and Kohlhagen, 1978; Auboin and Ruta, 2013). Lastly, we include the proportion of the population with internet access ($x_{it,\text{net}}^{(1)}/x_{jt,\text{net}}^{(2)}$),

Table A.18: Prediction RMSPE in five cities under the rolling window setting.

City	Model	Start Time Point				
		2010-Q1	2010-Q2	2010-Q3	2010-Q4	2021-Q1
Charlotte	GTNAR	0.0848	0.0841	0.0845	0.0848	0.0854
	sVAR	0.1145	0.1067	0.0999	0.1048	0.1061
	BiTR	0.0560	0.2329	0.7017	0.5040	0.6093
	GNAR-R	0.0964	0.1128	0.1227	0.0923	0.0949
	GNAR-C	2.6539	0.1023	0.2522	0.1080	0.0953
	LSTM	0.0866	0.0851	0.0863	0.0860	0.0871
LasVegas	GTNAR	0.0786	0.0760	0.0790	0.0814	0.0829
	sVAR	0.1118	0.1344	0.1300	0.1273	0.1442
	BiTR	1.5677	1.5457	1.6350	1.6465	1.2869
	GNAR-R	0.0862	0.1124	7.0599	0.6934	0.0849
	GNAR-C	43.9337	3.4841	37.4009	0.0943	0.2487
	LSTM	0.0789	0.0782	0.0790	0.0816	0.0837
Phoenix	GTNAR	0.0791	0.0762	0.0716	0.0690	0.0682
	sVAR	0.0877	0.0881	0.0861	0.0930	0.0881
	BiTR	0.6108	1.1864	0.2036	0.2322	0.2437
	GNAR-R	0.0801	0.0796	0.0792	0.0736	0.0733
	GNAR-C	1.3056	0.0795	0.1082	0.0745	0.1303
	LSTM	0.0782	0.0790	0.0738	0.0700	0.0717
Scottsdale	GTNAR	0.0617	0.0614	0.0611	0.0619	0.0600
	sVAR	0.0772	0.0720	0.0715	0.0766	0.0730
	BiTR	0.7890	2.0271	0.8091	0.8288	0.0446
	GNAR-R	0.0655	0.0639	0.0680	0.0638	0.0639
	GNAR-C	0.1423	0.0778	0.0656	0.0713	0.1265
	LSTM	0.0632	0.0624	0.0616	0.0624	0.0606
Toronto	GTNAR	0.0989	0.0972	0.0958	0.0915	0.0865
	sVAR	0.1117	0.1106	0.1105	0.1067	0.1083
	BiTR	1.3064	0.9224	0.6444	1.2483	0.9182
	GNAR-R	0.1013	0.1006	0.1013	0.1022	0.0951
	GNAR-C	5.1932	3.0433	6.2586	18.4431	0.4216
	LSTM	0.0999	0.0985	0.0969	0.0932	0.0877

a metric that captures the effects of the rapid development of e-commerce and its effect on cross-regional commodity circulation, as discussed in studies such as those by Freund and Weinhold (2004) and Terzi (2011). This covariate data has been sourced from both the World Bank Database¹ and the IMF-DOTS database, encompassing the years 1992 to 2021. Figure A.12 depicts the accumulated response for export volume in the year 2022 versus the three covariates of export countries in 2021, which suggests that a lower exchange rate, higher GDP, or a higher proportion of internet users in 2021 are associated with potentially higher export volumes in 2022.

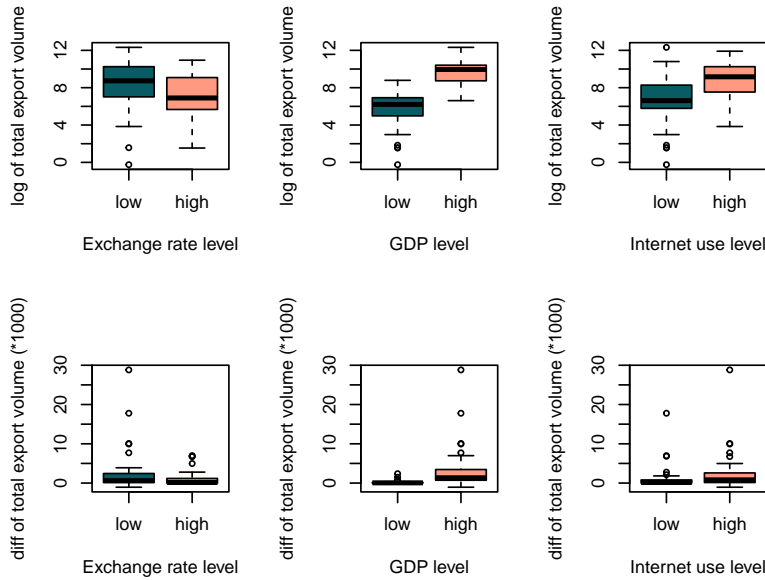


Figure A.12: Top panels: the relationship between the total export volume in 2022 and the levels of three covariates in 2021 (the categorization into "low" and "high" levels is determined by the medians). Bottom panels: the relationship between the first-order difference in exports between 2022 and 2021 and the levels of the same three covariates in 2021.

Spatial Networks. The countries under study are interconnected through various international relationships. In this analysis, our focus is on the spatial network spillover effects, a phenomenon well-documented in the literature (Redding and Venables, 2004; Yin et al., 2020). After filtering out countries/regions with missing data, we have a total of $N_1 = N_2 = 82$ observed units. To construct the adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$, we reference neighboring countries and territories for each country and region, as listed on Wikipedia²². In these adjacency matrices, an element $a_{ij}^{(1)} = 1$ ($a_{ij}^{(2)} = 1$) signifies the presence of at least one land border between the two regions. Consequently, our response variable, $Y_{ij,t}$, is indexed by two networks, the exporter network $\mathbf{A}^{(1)}$ and the importer network $\mathbf{A}^{(2)}$. In particular, we note that although these two networks share the same network nodes and

2. https://en.wikipedia.org/wiki/List_of_countries_and_territories_by_number_of_land_borders

topology, they exert different influences on the response variable, accounting for export and import spatial spillover effects on IEV, respectively.

N.2 Estimating Results

We next implement the GMNAR model to investigate the dynamic trading patterns around the world. The QIC suggests two groups for the exporter network ($\mathbf{A}^{(1)}$) nodes and three groups for the importer network ($\mathbf{A}^{(2)}$) nodes. The group patterns are color-coded in Figure A.13 and estimated model parameters are summarized in Table A.19. We summarize the estimated results as follows.

Self-momentum effects. Table A.19 reveals that the self-momentum effects $\hat{\alpha}$ are consistently negative for all cases, with the exception of $\hat{\alpha}_{22}$. One possible explanation for the negative auto-correlation is that the demand for exported goods from specific countries within the importer network tends to remain relatively stable over the years. Consequently, a significant increase in imports in the previous year may indicate a reduction in the importation of these goods in the following year. The sole exception is $\hat{\alpha}_{22} > 0$, representing the self-momentum of exported goods from group 2 of the exporter network (i.e., mainland China and Brazil) to group 2 of the importer network (i.e., Canada, mainland China and Tunisia). This implies that trade between these countries has consistently grown at an accelerated speed from 1992 to 2021. For instance, mainland China’s exports to the Canada have increased at annualized rates of 10.8% and 14.7%, respectively, while Brazil’s exports to mainland China and Canada have grown at rates of 17.6% and 8.41%³. Such observations align well with the estimation by the GTNAR method.

Spatial exporter network effects. Table A.19 indicates that the export network effects ($\hat{\lambda}_{g^{(1)}}^{(1)}$ ’s) are negative for both groups, but the second group, consisting of Brazil and mainland China, exhibits a substantially larger magnitude. The negative exporter network effects imply that a country’s export is negatively influenced when neighboring countries experience significant increases in their exports, particularly pronounced in the case of large developing countries like Brazil and mainland China. This observation suggests that countries in neighboring regions tend to be competitors rather than collaborators in the export market.

Spatial importer network effects. In Table A.19, we also observe that the import network effects $\hat{\lambda}_{g^{(2)}}^{(2)}$ are positive in the first two groups (i.e., $\hat{\lambda}_1^{(2)} = 0.042$ and $\hat{\lambda}_2^{(2)} = 0.032$), but negative in the third group (i.e., $\hat{\lambda}_3^{(2)} = -0.005$). The positive values of $\hat{\lambda}_{g^{(2)}}^{(2)}$ for the first two groups indicate that if neighboring countries import more goods, it positively influences their own imports. This implies a positive spillover effect of the import market for regions within the first two groups (such as mainland China, the United States and Canada). Conversely, the third group predominantly comprises developing countries in South America and Africa, which are shown to be less affected by their neighbors’ imports.

Covariates effects. In Table A.19, we can observe distinct effects on IEV (International Export Volume) resulting from the three covariates. Firstly, it is notable that for the second group of exporters, $\hat{\zeta}_{\text{exc},2}^{(1)} = -1.710$, signifying a significant negative correlation. This suggests that higher export costs lead to reduced IEV in Brazil and mainland China.

3. <https://oec.world/en/profile/bilateral-country>

Conversely, for the second group of importers, $\hat{\zeta}_{exc,2}^{(2)} = 0.035$, which is a positive value and aligns with expectations, as higher exchange rates generally correspond to increased import volumes. Secondly, both $\hat{\zeta}_{gdp}^{(1)}$ and $\hat{\zeta}_{gdp}^{(2)}$ display positive effects on IEV across all countries. This aligns with existing research indicating that domestic economic development bolsters international trade (Bernard and Jensen, 1999, 2004; Berman and Héricourt, 2010). Thirdly, we can observe that $\hat{\zeta}_{net,1}^{(1)} = 0.228$, suggesting that a higher level of Internet access in most exporting countries leads to increased IEV. However, it's worth noting that countries in the second export group (Brazil and mainland China), as well as the second importer group (mainland China and Canada), exhibit negative coefficients concerning Internet population percentage. This discrepancy may be attributed to the rapid growth of internet users in China, which has outpaced the growth in export volume from 1993 to 2022.

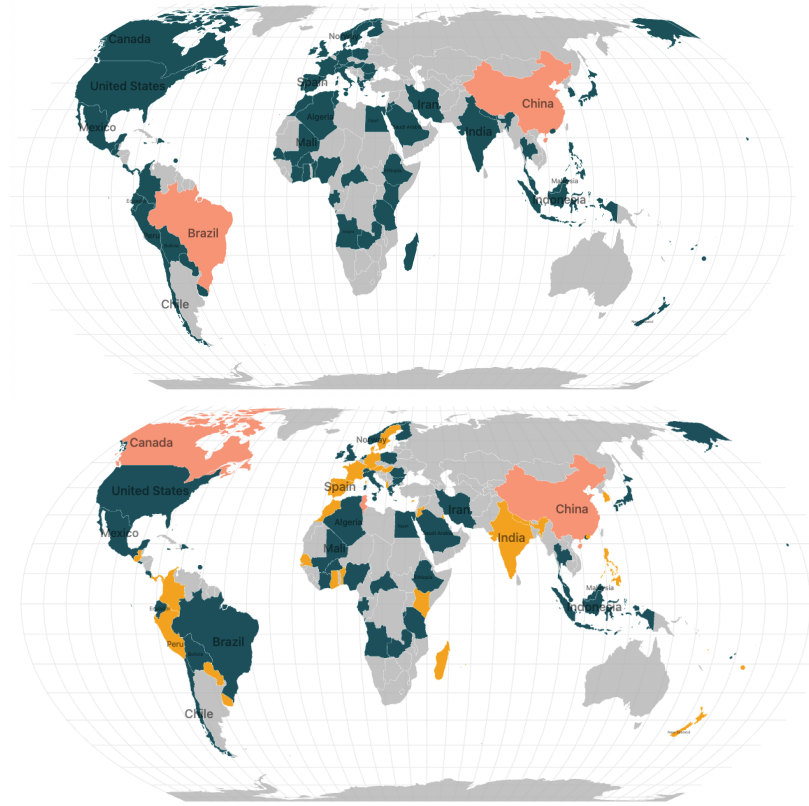


Figure A.13: The group patterns of export (the top) and import (the bottom) markets. The Group 1 and 2 in the export market are marked with dark green and pink. The Group 1, 2, 3 in the import market are marked with dark green, pink and orange, respectively. The gray color refers to the regions excluded in the model dataset.

Table A.19: Estimation results for trading data analysis. The p -values are shown in the parenthesis.

Parameters	$\lambda_{g^{(1)}}^{(1)}$		$\lambda_{g^{(2)}}^{(2)}$		
		-0.003 (<0.001)	-0.512 (<0.001)	0.042 (<0.001)	0.032 (<0.001)
	$\zeta_{g^{(1)}}^{(1)}$		$\zeta_{g^{(2)}}^{(2)}$		
Intercept	-0.011 (0.960)	2.072 (0.002)	-0.369 (0.112)	2.674 (<0.001)	-0.244 (0.329)
$\zeta_{\text{exc}}^{(1)}/\zeta_{\text{exc}}^{(2)}$	-0.120 (0.484)	-1.710 (0.036)	-0.104 (0.549)	0.035 (<0.001)	-1.399 (0.027)
$\zeta_{\text{gdp}}^{(1)}/\zeta_{\text{gdp}}^{(2)}$	0.693 (<0.001)	2.544 (<0.001)	1.046 (<0.001)	1.015 (<0.001)	1.681 (<0.001)
$\zeta_{\text{net}}^{(1)}/\zeta_{\text{net}}^{(2)}$	0.228 (0.022)	-2.962 (0.010)	0.029 (0.838)	-1.314 (0.001)	0.206 (0.186)
	$\alpha^T \in \mathbb{R}^{G_1 \times G_2}$				
		-0.431 (<0.001)		-0.494 (<0.001)	
		-0.287 (<0.001)		0.021 (<0.001)	
		-0.315 (<0.001)		-0.311 (<0.001)	