Learning cross-layer dependence structure in multilayer networks

Jiaheng Li¹, Jonathan Stewart¹

¹Department of Statistics, Florida State University, e-mail: jl20gx@fsu.edu; jrstewart@fsu.edu

Abstract: We propose a novel class of separable multilayer network models to capture cross-layer dependencies in multilayer networks, enabling the analysis of how interactions in one or more layers may influence interactions in other layers. Our approach separates the network formation process from the layer formation process, and is able to extend existing single-layer network models to multilayer network models that accommodate crosslayer dependence. We establish non-asymptotic and minimax-optimal error bounds for maximum likelihood estimators and demonstrate the convergence rate in scenarios of increasing parameter dimension. Additionally, we establish non-asymptotic error bounds for multivariate normal approximations and propose a model selection method that controls the false discovery rate. Simulation studies and an application to the Lazega lawyers network show that our framework and method perform well in realistic settings.

Keywords and phrases: Multilayer networks, statistical network analysis, social network analysis, network data, Markov random fields, graphical models.

1. Introduction

Multilayer networks have become a recent focal point of research in the field of statistical network analysis [e.g., 26, 9, 2, 24, 30, 11, 37, 19], arising in applications where a common set of elements in a population interact through multiple modes or relationships with other elements in the population. A prototypical example in the literature might be the Lazega law firm network [25], in which attorneys are linked through various forms of interaction, such as advice seeking, friendship, collaboration, etc., each of which would form a distinct layer in the multilayer network [24]. In essence, a multilayer network is a composite structure, where each layer captures a specific type of interaction or relationship between the same set of elements.

Edges in one layer of the multilayer network may depend on edges in other layers, creating what is known as cross-layer dependence. Understanding the drivers of edge formation in multilayer networks requires learning the dependence structures across these layers. A key challenge lies in the fact that the cross-layer dependence can be highly varied and complex, and the development of statistical models with theoretical guarantees for network data with dependent edges is challenging. Current methodological frameworks for multilayer networks can be broadly categorized into two main groups:

- 1. Statistical models equipped with theoretical guarantees often rely on latent variable constructions [e.g., 30, 2, 19]. These models typically assume conditional independence of edges given the latent variables, following standard practices within the field.
- Statistical models that do not provide formal theoretical guarantees [e.g., 9, 24]. Instead, these methods extend existing approaches by explicitly allowing for edge dependencies, thereby relaxing the conditional independence assumptions present in the first class of models.

In this work, we address a critical gap in the literature by introducing a separable multilayer network modeling framework for multilayer networks. Our approach not only accommodates dependent edges but also provides theoretical guarantees for both estimation and inference without relying on any latent variables. Specifically, we extend single-layer network models to the multilayer setting, with a central focus on identifying and understanding cross-laver dependence structures. A key advantage of our proposed framework is that we are able to distinguish the network formation process from the layer formation process. This allows us to create a wide range of novel multilayer network models derived from established single-layer network models, such as exponential-family random graph models, stochastic block models, and latent space models. By employing Markov random field specifications, we develop adaptable and comprehensive models to capture cross-layer dependencies in multilayer networks. As a result, our framework jointly models both network structures and crosslayer dependence, thus enabling any single-layer network model to be extended to the multilayer setting. Our main contributions in this work include:

- 1. Introducing a novel framework for modeling cross-layer dependence in multilayer networks that synchronizes with current network models in the literature.
- 2. Deriving non-asymptotic theoretical guarantees in scenarios where the number of parameters tends to infinity, which establishes bounds on the:
 - (a) Statistical error of maximum likelihood estimators.
 - (b) Error of the multivariate normal approximation of estimators.
- 3. Elaborating a model selection algorithm which controls the false discovery rate.

The rest of the paper is organized as follows. Section 2 introduces our modeling framework and includes illustrative examples. The consistency and the minimax optimal results are contained in Section 3. The multivariate normal approximation theory is presented in Section 4. The results of simulation studies are provided in Section 5, together with different testing procedures for model selection which control the false discovery rate. An application of our developed framework and methodology is given in Section 6, concluding with a discussion presented in Section 7. The code and data to reproduce the simulations and analyses can be found in our package online.¹

¹https://github.com/jiaheng-li/mlyrnetwork

2. Modeling cross-layer dependence in multilayer networks

A multilayer network can be represented as a sequence of $1 \leq K < \infty$ random graphs $\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(K)}$ each defined on a common set of $N \geq 3$ nodes, which we take without loss to be the set $\mathcal{N} = \{1, \ldots, N\}$. We call the graphs $\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(K)}$ the *layers* of the network, and represent the multilayer network as the quantity $\mathbf{X} = (\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(K)})$.

Connections between pairs of nodes $\{i, j\} \subset \mathbb{N}$ in each layer $k \in \{1, \dots, K\}$ are modeled by random variables

$$X_{i,j}^{(k)} = \begin{cases} 1 & \text{nodes } i \text{ and } j \text{ are connected in layer } k \\ 0 & \text{otherwise} \end{cases}$$

We refer to all connections of a pair of nodes $\{i, j\} \subset \mathbb{N}$ across the K layers as a *dyad* which we denote by $\mathbf{X}_{i,j} = (X_{i,j}^{(1)}, \ldots, X_{i,j}^{(K)}) \in \{0,1\}^K$. A multilayer network can be represented by a collection of dyads as $\mathbf{X} = (\mathbf{X}_{i,j})_{\{i,j\} \subset \mathbb{N}}$ alternatively.

For notational ease, we will consider undirected multilayer networks, which imply that the network layers $\mathbf{X}^{(1)}, \ldots, \mathbf{X}^{(K)}$ are undirected random graphs; extensions to directed multilayer networks or mixed multilayer networks with both directed and undirected layers will typically be straightforward, involving only notational adaptations in subscripts in most cases. We adopt the usual conventions for undirected networks, i.e., we assume that $X_{i,j}^{(k)} = X_{j,i}^{(k)}$ (all $\{i,j\} \subset \mathbb{N}, 1 \leq k \leq K$) and $X_{i,i}^{(k)} = 0$ (all $i \in \mathbb{N}, 1 \leq k \leq K$). The sample space of each layer $\mathbf{X}^{(k)}$ is therefore the product space $\mathbb{X}^{(k)} \coloneqq \{0,1\}^{\binom{N}{2}}$ (k = $1, \ldots, K$), and the sample space \mathbb{X} of \mathbf{X} is the product space of the sample spaces of the individual layers, i.e., $\mathbb{X} \coloneqq \mathbb{X}^{(1)} \times \cdots \times \mathbb{X}^{(K)}$. The sample space of dyad $\{i, j\} \subset \mathbb{N}$ is the product space $\mathbb{X}_{i,j} \coloneqq \{0,1\}^K$.

A challenge in the statistical modeling of network data lies in the fact that networks have many distinguishing properties, including:

- 1. **Sparsity.** Many real-world networks are sparse, in the sense that the expected number of edges in the network grows at a rate slower than $\binom{N}{2}$. The phenomena of network sparsity manifests in a variety of different applications, usually due to constraints, such as time or financial constraints, which can limit the number of connections any node can maintain at a given point in time [22, 8].
- 2. Node heterogeneity. Different actors in a social network will have different properties, called node covariates, which can lead to different propensities to form edges. A key example is assortative and disassortative mixing patterns in networks [31, 23], as well as differences in structural patterns in the network [1, 27].
- 3. Edge dependence. In addition to node-based effects that give rise to heterogeneity in propensities for nodes to form edges, scientific and statistical evidence suggests edges are dependent in many applications [18, 13, 6],

and modeling a single system of multiple binary random variables without replication is a challenging statistical problem inherent to many statistical network analysis applications.

Each of the above gives rise to distinct challenges for modeling network data and performing statistical inference in statistical network analysis applications, and it is not straightforward to construct models that due justice to each of these and more. To address these challenges, a plethora of statistical models have been proposed to model network data, which for single-layer networks have included exponential-families of random graph models [e.g., 28, 36], stochastic block models [e.g., 17], latent metric space models [e.g., 16], random dot product graphs [e.g., 3], exchangeable random graph models [e.g., 10, 12], and more. In this work, we build upon the many classes of single-layer network data models by introducing a separable multilayer network modeling framework. This framework enables existing single-layer network models to be extended to the multilayer setting and simultaneously enables learning cross-layer dependence and interactions across different layers in the multilayer network.

2.1. Separable multilayer network models

Multilayer networks are subject to the same forces and phenomena as single layer networks, as multiple modes of relation or interaction do not remove constraints or properties of nodes which are fundamental to network data applications. The same set of nodes is defined across all layers in a multilayer network, and because all layers share the same set of nodes, the dyadic connections among these nodes fundamentally define the network formation process. By specifying a single-layer network as the foundational structure reference, we can separate the network formation process from the layer formation process. In doing so, the single-layer network serves as the baseline for establishing dyadic relationships that represent the relational structure across all layers of the multilayer network. As a result, the network formation process determines which dyads have the potential to form connections, i.e., which pair of nodes may exhibit at least one edge in any of the layers. In contrast, The layer formation process dictates the particular layers in which these connections appear. To learn the effects of cross-layer dependence in multilayer networks, we propose the class of separable multilayer network models, which extend the broad literature on single-layer network models into the multilayer realm. These models can incorporate an arbitrary single-layer network structure as the foundational baseline and ensures that the underlying single-layer network can be recovered from observations of the multilayer network. We illustrate this approach and its advantages through our proposed modeling framework. We specify probability distributions on a double of networks (X, Y), where Y will represent the network formation process, which we will call the *basis network*, and X will represent the realized multilayer network. We assume that $\boldsymbol{Y} \in \mathbb{Y} \coloneqq \{0,1\}^{\binom{N}{2}}$ is an undirected, single-layer network defined on the set of nodes \mathcal{N} where, for

all
$$\{i, j\} \subset \mathbb{N}$$
,
 $Y_{i,j} = \begin{cases} 1 & \text{nodes } i \text{ and } j \text{ are connected in the basis network} \\ 0 & \text{otherwise} \end{cases}$,

making the usual conventions for undirected networks mentioned previously. For (\mathbf{X}, \mathbf{Y}) , we consider semi-parametric families of probability distributions $\mathcal{F} := \{\mathbb{P}_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \mathbb{R}^p\}$ which are absolutely continuous with respect to a σ -finite measure ν defined on $\mathcal{P}(\mathbb{X} \times \mathbb{Y})$, where $\mathcal{P}(\mathbb{X} \times \mathbb{Y})$ is the power set of $\mathbb{X} \times \mathbb{Y}$. Typically, ν will be the counting measure. We say the probability mass function $\mathbb{P}_{\boldsymbol{\theta}} \in \mathcal{F}$ defines a *separable multilayer network model* if $\mathbb{P}_{\boldsymbol{\theta}}$ admits the form:

$$\mathbb{P}_{\boldsymbol{\theta}}(\{(\boldsymbol{x},\boldsymbol{y})\}) = f(\boldsymbol{x},\boldsymbol{\theta}) g(\boldsymbol{y}) h(\boldsymbol{x},\boldsymbol{y}) \psi(\boldsymbol{\theta},\boldsymbol{y}), \qquad (\boldsymbol{x},\boldsymbol{y}) \in \mathbb{X} \times \mathbb{Y}, \qquad (1)$$

where

• $f: \mathbb{X} \times \mathbb{R}^p \mapsto (0, 1)$ is given by

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \prod_{\{i,j\} \subset \mathcal{N}} \exp \left(\sum_{k=1}^{K} \theta_k x_{i,j}^{(k)} + \sum_{k$$

where $H \leq K$ is the highest order of cross-layer interactions included in the model. We write $\theta_{k_1,k_2,\ldots,k_h}$ to reference the *h*-order interaction parameter for the interaction term among layers $\{k_1,\ldots,k_h\} \subseteq \{1,\ldots,K\}$.

- $g: \mathbb{Y} \mapsto (0,1)$ is the marginal probability mass function of Y and is assumed to be strictly positive on \mathbb{Y} .
- $h: \mathbb{X} \times \mathbb{Y} \mapsto \{0, 1\}$ is given by

$$h(\boldsymbol{x}, \boldsymbol{y}) = \prod_{\{i,j\} \subset \mathbb{N}} \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0)^{y_{i,j}} \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 = 0)^{1-y_{i,j}},$$

where $\boldsymbol{x}_{i,j} = (x_{i,j}^{(1)}, \dots, x_{i,j}^{(K)}) \in \mathbb{X}_{i,j} \ (\{i, j\} \subset \mathbb{N}).$

• $\psi: \Theta \times \mathbb{Y} \mapsto (0, \infty)$ is defined by

$$\psi(\boldsymbol{\theta}, \boldsymbol{y}) = \left[\sum_{\boldsymbol{x} \in \mathbb{X}} f(\boldsymbol{x}, \boldsymbol{\theta}) h(\boldsymbol{x}, \boldsymbol{y})\right]^{-1},$$

ensuring (1) will be a valid probability mass function for (X, Y).

The notation $\mathbb{P}_{\theta}(\{(x, y)\})$ is well-defined for each pair $(x, y) \in \mathbb{X} \times \mathbb{Y}$, as \mathbb{P}_{θ} is a probability measure defined on $\mathcal{P}(\mathbb{X} \times \mathbb{Y})$. Frequently, we will write the



Fig 1: Multilayer networks specified by three different basis network structures: the latent space model (LSM), the exponential random graph model (ERGM), and the stochastic block model (SBM).

probability expressions $\mathbb{P}_{\theta}(X = x, Y = y)$ for the joint probability of $\{(x, y)\}$, and $\mathbb{P}_{\theta}(X = x | Y = y)$ for the conditional probability of the event X = xconditional on the event Y = y. We denote the data-generating parameter vector by $\theta^* \in \mathbb{R}^p$, and the corresponding probability measure and expectation operator by $\mathbb{P} \equiv \mathbb{P}_{\theta^*}$ and $\mathbb{E} \equiv \mathbb{E}_{\theta^*}$, respectively.

The specification in equation (1) separates the network formation process **Y**, specified by $q(\mathbf{y})$, from the layer formation process, specified by $f(\mathbf{x}, \boldsymbol{\theta})$. The two are joined by the function $h(\boldsymbol{x}, \boldsymbol{y})$, which ensures $\|\boldsymbol{x}_{i,j}\|_1 = 0$ whenever $Y_{i,j} = 0$ and $\|\boldsymbol{x}_{i,j}\|_1 > 0$ whenever $Y_{i,j} = 1$, as we allow edges between nodes $i \in \mathbb{N}$ and $j \in \mathbb{N}$ in \boldsymbol{X} if and only if $Y_{i,j} = 1$. We call dyads $\{i, j\} \subset \mathbb{N}$ with $Y_{i,j} = 1$ activated dyads, and a pair $(\boldsymbol{x}, \boldsymbol{y}) \in \mathbb{X} \times \mathbb{Y}$ that satisfies $h(\boldsymbol{x}, \boldsymbol{y}) = 1$ is said to be a *concordant* pair. We will only focus on concordant pairs of multilayer networks since our modeling framework guarantees the recovery of the basis network Y given an observation of X, a point that will be made clear shortly in Proposition 1. The function $\psi(\theta, y)$ ensures the resulting product of functions will be a valid probability mass function, and it has less of a direct role in modeling the cross-layer dependence, essentially fulfilling the role of a normalizing constant for the conditional probability distribution of X given Y, as derived in Proposition 1. Such specifications have the advantage of being able to specify the network formation process separately from the process that populates the layers of activated dyads, thus modeling the cross-layer dependence conditional on Y. To illustrate the flexibility and generality of (1), observe that $q(\mathbf{y})$ is allowed to be any probability mass function for a single-layer network Y (e.g., exponential-family random graph model, stochastic block model, latent space model), provided $q(\mathbf{y}) > 0$ for all $\mathbf{y} \in \mathbb{Y}$. To illustrate this point, Figure

6

1 displays various multilayer networks with K = 3 layers where the basis network is specified via three different models, demonstrating that our modeling framework is capable of ensuring that the multilayer network respects structural properties of the underlying basis network. We view our framework as semiparametric as $g(\mathbf{y})$ need not assume a specific parametric form. Moreover, our framework can be viewed as non-parametric when the maximal possible order of interaction terms are included in (1), a point on which we further elaborate later. An important feature of our framework lies in the fact that the choice of the probability distribution for the network formation process does not directly influence the estimation for the cross-layer dependence structure, i.e., the choice of $g(\mathbf{y})$ does not directly influence estimation for $\boldsymbol{\theta}^*$. Proposition 1 demonstrates this point in the case of likelihood-based inference.

Proposition 1 Let $\{\mathbb{P}_{\theta} : \theta \in \mathbb{R}^p\}$ satisfy (1). Then the following hold:

- 1. For each $x \in \mathbb{X}$, Y = y (\mathbb{P}_{θ} -a.s.) for one and only one $y \in \mathbb{Y}$.
- 2. **Y** is predictable via **X**, i.e., for each $x \in X$, $\mathbb{P}_{\theta}(Y = y | X = x) = 1$ where

$$y_{i,j} = \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0), \quad \{i,j\} \subset \mathcal{N}.$$

3. For all $(\boldsymbol{x}, \boldsymbol{y}) \in \mathbb{X} \times \mathbb{Y}$ with $h(\boldsymbol{x}, \boldsymbol{y}) = 1$,

$$\log \mathbb{P}_{\boldsymbol{ heta}}(\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Y} = \boldsymbol{y}) = \log \mathbb{P}_{\boldsymbol{ heta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) + \log g(\boldsymbol{y}),$$

where $\mathbb{P}_{\theta}(X = x | Y = y)$ belongs to a minimal exponential family with natural parameter vector $\theta \in \mathbb{R}^p$ and is given by

$$\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) = \exp(\log f(\boldsymbol{x}, \boldsymbol{\theta}) + \log \psi(\boldsymbol{\theta}, \boldsymbol{y})).$$

Proposition 1 establishes a few key facts for the inference of cross-layer dependence structures in multilaver networks. First, we are able to observe Ythrough X, as given any observation $x \in \mathbb{X}$ of the multilayer network X, $\mathbb{P}_{\theta}(Y = y \mid X = x) = 1$ for one, and only one, $y \in \mathbb{Y}$. In other words, through the observation of \boldsymbol{x} , we can infer with probability 1 the corresponding \boldsymbol{y} due to the form of (1). The significance of this result is that we do not need to treat the basis network \boldsymbol{Y} as a latent network, which would require additional statistical and computational methodology to handle the latent missing network data. Second, we see that the inference for θ^* is unaffected by the choice of g(y); although, the statistical guarantees for estimators of θ^{\star} will be indirectly influenced by the choice of $q(\mathbf{y})$, a point which we discuss in later sections. Moreover, the above choice for $f(x, \theta)$ and the functional form of $\mathbb{P}_{\theta}(X = x \mid Y = y)$ derived in Proposition 1 establishes that $\log \mathbb{P}_{\theta}(X = x | Y = y)$ corresponds to the log-likelihood of a minimal exponential family, accessing a broad literature of statistical methodology and theory [e.g., 41]. We note that other specifications for $f(\boldsymbol{x}, \boldsymbol{\theta})$ are possible, but that Markov random field specifications provide a powerful class of models for dependent data [e.g., 44], and in the case of the saturated model with maximal interaction term H = K, it completely specifies all possible probabilities of outcomes $x_{i,j} \in \{0,1\}^K$, presenting a non-parametric model class for multilayer networks.

2.2. Example of a multilayer network with pairwise interactions

We illustrate cross-layer dependence among layers in our modeling framework by considering a separable multilayer network model using the Markov random field specification for $f(\boldsymbol{x}, \boldsymbol{\theta})$ given in the previous section and maximal interaction term H = 2:

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \prod_{\{i,j\} \subset \mathcal{N}} \exp\left(\sum_{k=1}^{K} \theta_k x_{i,j}^{(k)} + \sum_{k(2)$$

The dimension of the parameter vector $\boldsymbol{\theta}$ is dim $(\boldsymbol{\theta}) = K + \binom{K}{2}$, with K parameters governing the single-layer effects for the K layers and $\binom{K}{2}$ combinations of layers to form the pairwise interactions for the cross-layer dependence effects. Define the (K-1)-dimensional vector $X_{i,j}^{(-k)} \coloneqq (X_{i,j}^{(l)}) : l \in \{1, \ldots, K\} \setminus \{k\}$ to

Define the (K-1)-dimensional vector $X_{i,j}^{(-k)} := (X_{i,j}^{(l)} : l \in \{1, \ldots, K\} \setminus \{k\})$ to be the vector of edge variables in $X_{i,j}$ which excludes the edge variable $X_{i,j}^{(k)}$, i.e., excluding the edge variable between nodes i and j in layer k. The conditional log-odds of edge $X_{i,j}^{(k)}$ takes the form:

$$\log \frac{\mathbb{P}(X_{i,j}^{(k)} = 1 \mid \boldsymbol{X}_{i,j}^{(-k)} = \boldsymbol{x}_{i,j}^{(-k)}, Y_{i,j} = 1)}{\mathbb{P}(X_{i,j}^{(k)} = 0 \mid \boldsymbol{X}_{i,j}^{(-k)} = \boldsymbol{x}_{i,j}^{(-k)}, Y_{i,j} = 1)} = \begin{cases} \theta_k + \sum_{l \neq k}^K \theta_{k,l} x_{i,j}^{(l)}, & \|\boldsymbol{x}_{i,j}^{(-k)}\|_1 > 0 \\ +\infty, & \|\boldsymbol{x}_{i,j}^{(-k)}\|_1 = 0 \end{cases}$$

A primary advantage and motivation of using a parametric Markov random field specification for $f(\boldsymbol{x}, \boldsymbol{\theta})$ lies in the interpretability of the model. An effective approach to analyzing and understanding marginal network effects in such specifications is to study conditional log-odds of edges under different conditioning statements [e.g., 39]. By the form of $h(\boldsymbol{x}, \boldsymbol{y})$, when $Y_{i,j} = 1$, we require $\|\boldsymbol{x}_{i,j}\|_1 > 0$, meaning nodes i and j must have at least one connection in \boldsymbol{X} . This is seen through the log-odds formula above, where the log-odds of edge $X_{i,j}^{(k)}$ is equal to $+\infty$ when $\|\boldsymbol{x}_{i,j}^{(-k)}\|_1 = 0$. In contrast, when $\|\boldsymbol{x}_{i,j}^{(-k)}\|_1 > 0$, the constraint $\|\boldsymbol{x}_{i,j}\|_1 > 0$ is already satisfied, and the log-odds of edge $X_{i,j}^{(k)}$ depends on the layer specific parameter θ_k , as well as the pairwise interaction effects where edges present in other layers $l \in \{1, \ldots, K\} \setminus \{k\}$ can influence the likelihood of the edge $X_{i,j}^{(k)}$ depending on the signs and magnitudes of the pairwise interaction parameters $\theta_{k,l}$ ($\{k, l\} \subseteq \{1, \ldots, K\}$).

3. Estimation of cross-layer dependence structure

For separable multilayer network models satisfying (1), Proposition 1 establishes that the log-likelihood function takes the form

$$\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \coloneqq \log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Y} = \boldsymbol{y}) \\ = \log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) + \log g(\boldsymbol{y}).$$
(3)

Given an observation $x \in \mathbb{X}$ of the multilayer network X, and therefore an observation $y \in \mathbb{Y}$ of Y by Proposition 1, we denote the set of maximum likelihood estimators by

$$\widehat{oldsymbol{\Theta}} \hspace{.1in}\coloneqq \hspace{.1in} \left\{ oldsymbol{ heta} \in \mathbb{R}^p \, : \, \ell(oldsymbol{ heta}; oldsymbol{x}, oldsymbol{y}) = \sup_{oldsymbol{ heta}' \in \mathbb{R}^p} \, \ell(oldsymbol{ heta}'; oldsymbol{x}, oldsymbol{y})
ight\},$$

and reference individual elements of the set by $\hat{\theta} \in \hat{\Theta}$. As Proposition 1 establishes log $\mathbb{P}_{\theta}(X = x | Y = y)$ to be a minimal, and by construction regular, exponential family, $|\hat{\Theta}| \in \{0, 1\}$, i.e., when the maximum likelihood estimator exists, the set $\hat{\Theta}$ will contain a unique element when non-empty [Proposition 3.11, pp. 32–33, 41]. As seen from the forms of $\ell(\theta; x, y)$ given above, the gradients and Hessians of the log-likelihood equations do not directly depend on g(y). However, the following lemma shows how theoretical guarantees for estimators of θ^* will be indirectly influenced by the choice of g(y).

Lemma 1. Consider a family $\{\mathbb{P}_{\theta} : \theta \in \mathbb{R}^p\}$ of separable multilayer network models satisfying (1) and an observation $x \in \mathbb{X}$ of X. Let (x, y) be the concordant pair where y is given by Proposition 1. Define, for each pair of nodes $\{i, j\} \subset \mathbb{N}$,

$$L_{i,j}(\boldsymbol{\theta}, \boldsymbol{x}_{i,j}, \boldsymbol{y}) \cong \log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X}_{i,j} = \boldsymbol{x}_{i,j} \mid \boldsymbol{Y} = \boldsymbol{y}).$$

Then there exists a $p \times p$ matrix $I(\boldsymbol{\theta})$ such that

$$\mathbb{E}\left[-\nabla_{\boldsymbol{\theta}}^{2}L_{i,j}(\boldsymbol{\theta},\boldsymbol{X}_{i,j},\boldsymbol{Y}) \,|\, \boldsymbol{Y} = \boldsymbol{y}\right] = \begin{cases} I(\boldsymbol{\theta}) & Y_{i,j} = 1\\ \mathbf{0}_{p,p} & Y_{i,j} = 0, \end{cases}$$

for all $\{i, j\} \subset \mathbb{N}$, where $\mathbf{0}_{p,p}$ is the $p \times p$ matrix with all 0 entries, and

$$\begin{split} \lambda_{\min}(-\mathbb{E}\,\nabla^2_{\boldsymbol{\theta}}\,\ell(\boldsymbol{\theta};\boldsymbol{X},\boldsymbol{Y})) &= \lambda_{\min}(I(\boldsymbol{\theta}))\,\mathbb{E}\,\|\boldsymbol{Y}\|_1\\ \lambda_{\max}(-\mathbb{E}\,\nabla^2_{\boldsymbol{\theta}}\,\ell(\boldsymbol{\theta};\boldsymbol{X},\boldsymbol{Y})) &= \lambda_{\max}(I(\boldsymbol{\theta}))\,\mathbb{E}\,\|\boldsymbol{Y}\|_1, \end{split}$$

where $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$ are the smallest and the largest eigenvalue of matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$, respectively.

In classical settings with independent and identically distributed observations, the expected negative Hessian of the log-likelihood function is the Fisher information matrix and is expected to scale with the number of observations. In such cases, standard matrix theory indicates that the smallest eigenvalue of this expected negative Hessian matrix will scale with the sample size, provided the smallest eigenvalue of the Fisher information matrix is bounded from below. Lemma 1 extends this notion by establishing similar scaling behavior concerning the expected number of activated dyads $\mathbb{E} \| \boldsymbol{Y} \|_1$, proxying as an effective sample size. Analogously, $I(\boldsymbol{\theta})$ can be seen as the Fisher information of the population distribution governing individual activated dyads in \boldsymbol{Y} , mirroring the role of Fisher information for population distributions in classical independent and identically distributed scenarios.

Before we present our theoretical guarantees for maximum likelihood estimators in Theorem 1, we define some notations and outline some regularity assumptions for our theorem to follow. As we will show in Theorem 1, the choice of g(y) influences the estimation error through the expected number of edges in Y and through the covariances of edge variables in Y. Define

$$D_g := \sum_{\{i,j\} \prec \{v,w\} \subset \mathcal{N}} \mathbb{C}(Y_{i,j}, Y_{v,w}),$$

where $\{i, j\} \prec \{v, w\}$ implies the sum is taken with respect to the lexicographical ordering of pairs of nodes. Define $[D_g]^+ := \max\{0, D_g\}$ to be the positive part of D_g . Let $\epsilon > 0$ be fixed independent of N and p, and denote the ϵ -ball of the data-generating parameter θ^* by $\mathcal{B}_2(\theta^*, \epsilon) = \{\theta \in \mathbb{R}^p : \|\theta^* - \theta\|_2 \le \epsilon\}$. Define

$$\widetilde{\lambda}_{\min}^{\epsilon} \coloneqq \inf_{\boldsymbol{\theta} \in \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \epsilon)} \lambda_{\min}(I(\boldsymbol{\theta})) \quad \text{and} \quad \widetilde{\lambda}_{\max}^{\star} \coloneqq \lambda_{\max}(I(\boldsymbol{\theta}^{\star})),$$

where $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$ are the smallest and the largest eigenvalue of matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$, respectively.

Assumption 1. Assume there exists a $C_0 > 0$ such that $\mathbb{E} \| \mathbf{Y} \|_1 \ge 1$ and

$$\frac{[D_g]^+}{\mathbb{E} \|\boldsymbol{Y}\|_1} \leq C_0,$$

for all network sizes N.

Assumption 2. Assume the parameter dimension p satisfies

$$p \leq \sqrt{\widetilde{\lambda}_{\max}^{\star} \mathbb{E} \| \boldsymbol{Y} \|_{1}},$$

for all network sizes N.

Assumption 3. Assume that $\tilde{\lambda}_{\max}^{\star}$ and $\tilde{\lambda}_{\min}^{\epsilon}$ satisfy, as a function of the network size N,

$$\frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} = o\left(\sqrt{\frac{\mathbb{E}\|\boldsymbol{Y}\|_{1}}{p}}\right).$$

Assumptions 1–3 provide a foundation for Theorem 1 to establish the consistency result of the maximum likelihood estimator in large network settings. Assumption 1 imposes a lower bound on the expected number of activated dyads relative to the covariance as the network size N grows. Assumption 2 restricts the growth rate of p in relation to the network size and the largest eigenvalue of the Fisher information $I(\theta^*)$. Finally, Assumption 3 sets a constraint on the ratio between $\sqrt{\tilde{\lambda}_{\max}^*}$ and $\tilde{\lambda}_{\min}^{\epsilon}$, balancing eigenvalue magnitudes in a way that preserves estimator consistency under increasing network size. **Theorem 1.** Consider a multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ nodes. If Assumptions 1, 2, and 3 are satisfied, there exists $N_0 \ge 3$ such that, for all $N \ge N_0$, with probability at least $1 - \exp(-2p) - (\mathbb{E} \|\mathbf{Y}\|_1)^{-1}$, the set $\widehat{\mathbf{\Theta}}$ is non-empty and the unique element $\widehat{\boldsymbol{\theta}} \in \widehat{\mathbf{\Theta}}$ satisfies

$$\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2} \leq C \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{p}{\mathbb{E}\|\boldsymbol{Y}\|_{1}}}, \qquad (4)$$

11

where C > 0 is a constant independent of N and p.

The results of Theorem 1 establish a few key facts concerning statistical estimation of the data-generating parameter vector θ^{\star} . First, we can view the quantity $\tilde{\lambda}_{\min}^{\epsilon} \sqrt{\mathbb{E} \| \boldsymbol{Y} \|_1} / \tilde{\lambda}_{\max}^{\star}$ as the effective sample size in order to compare our results to classical settings with independent and identically distributed data. The effective sample size, together with the dimension of the model p, helps to determine the rate of convergence (with respect to the Euclidean distance) of maximum likelihood estimators. As previously mentioned, the quantity $\mathbb{E} \| \boldsymbol{Y} \|_{1}$ is determined by properties of $q(\mathbf{y})$, the marginal probability mass function of **Y**. While the specification of $q(\mathbf{y})$ does not directly influence the estimation algorithm, the statistical guarantees of estimators will depend on $q(\mathbf{u})$ producing enough activated dyads and not possessing overly strong dependence among edges in the single-layer basis network Y (Assumption 1). The requirement (Assumption 3) that the right-hand side of the bounds in Theorem 1 tends to 0 as $N \to \infty$ ensures that all regularity assumptions remain valid. Namely, key to our approach lies in the ability to control minimum eigenvalues of matrices $I(\boldsymbol{\theta})$ in a neighborhood of the data-generating parameter vector $\boldsymbol{\theta}^{\star}$. The condition that the bounds tend to 0 ensures that it is sufficient to control the smallest eigenvalue in a bounded set, i.e., we may let ϵ be fixed independent of N and p, and moreover, to ensure consistency in the sense that $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2} \to 0$ with probability approaching 1 as $N, p \to \infty$.

Corollary 1. Under the assumptions of Theorem 1, and in the case that the parameter dimension p is fixed, there exists $N_0 \geq 3$ such that, for all $N \geq N_0$ and $\alpha_N \in (2(\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}, 1/2)$, with probability at least $1 - \alpha_N$, the set $\widehat{\boldsymbol{\Theta}}$ is non-empty and the unique element $\widehat{\boldsymbol{\theta}} \in \widehat{\boldsymbol{\Theta}}$ satisfies

$$\|\widehat{\boldsymbol{ heta}} - \boldsymbol{ heta}^{\star}\|_{2} \leq C \left|\log\left(\frac{lpha_{N}}{2}\right)\right| \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_{1}}},$$

where C > 0 is a constant independent of N.

The corollary builds on the consistency result established in Theorem 1 by providing a similar bound in the situation that the parameter dimension p remains fixed. This simplifies the convergence rate by removing the dependence of p in the error term, which can yield sharper asymptotic guarantees. The introduction of the probability bound α_N offers an explicit control over the confidence level for the estimate's accuracy, which improves interpretability and practical applicability in finite samples. The bound on $\|\widehat{\theta} - \theta^{\star}\|_2$ in Corollary 1 now depends logarithmically on α_N , introducing a trade-off between the confidence level and the convergence rate. While the key dependencies remain on the effective sample size $\widetilde{\lambda}_{\min}^{\epsilon} \sqrt{\mathbb{E} \|\boldsymbol{Y}\|_1 / \widetilde{\lambda}_{\max}^{\star}}$ as in Theorem 1, Corollary 1 provides a useful refinement of the consistency result when the model's dimensionality is constrained.

We next present that the upper bound in Theorem 1 is minimax optimal up to a constant. Define the minimax risk to be

$$\mathcal{R}_N \coloneqq \inf_{\widehat{\boldsymbol{ heta}}} \sup_{\boldsymbol{ heta} \in \mathbb{R}^p} \mathbb{E}_{\boldsymbol{ heta}} \| \widehat{\boldsymbol{ heta}} - \boldsymbol{ heta} \|_2,$$

and

$$\widetilde{\lambda}_{\max}^{\epsilon} \coloneqq \sup_{\boldsymbol{\theta} \in \mathfrak{B}_{2}(\boldsymbol{\theta}^{\star}, \epsilon)} \lambda_{\max}\left(I(\boldsymbol{\theta})\right),$$

where $\lambda_{\max}(\mathbf{A})$ is the largest eigenvalue of matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$.

Theorem 2. Consider a separable multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ nodes. If Assumptions 1, 2 and 3 are satisfied, there exists a constant C > 0 independent of N and p, such that, the lower bound of the minimax risk \mathcal{R}_N satisfies

$$\mathcal{R}_N \geq C \frac{\widetilde{\lambda}_{\min}^{\epsilon}}{\widetilde{\lambda}_{\max}^{\epsilon}} \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{p}{\mathbb{E} \|\boldsymbol{Y}\|_1}}.$$

Theorem 2 establishes a lower bound for the minimax risk \mathcal{R}_N , differing from the upper bound of the ℓ_2 -error for the maximum likelihood estimator in Theorem 1 by a factor of $\tilde{\lambda}_{\min}^{\epsilon} / \tilde{\lambda}_{\max}^{\epsilon}$. Building on this result, we establish conditions for the minimax optimality of the maximum likelihood estimators in Corollary 2.

Corollary 2. Under the assumptions of Theorem 2 and the assumption that

$$\widetilde{\lambda}_{\max}^{\epsilon} = O\left(\widetilde{\lambda}_{\min}^{\epsilon}\right),\tag{5}$$

the maximum likelihood estimator $\hat{\theta}$ achieves the minimax rate of convergence, in the sense that the upper bound on the ℓ_2 -error of $\hat{\theta}$ given in Theorem 1 matches the lower bound of the minimax risk \mathcal{R}_N in Theorem 2, up to a constant.

The condition in (5) ensures that the rate of convergence for the maximum likelihood estimator achieves the minimax optimality by imposing a more direct and stringent relationship between the minimum and maximum eigenvalues than that required by Assumption 3. The control on the minimum and maximum eigenvalues for high-dimensional graphical models are common [e.g., 35, 45, 34], ensuring that the minimum and maximum eigenvalues of the information matrices within a neighborhood of the data-generating parameter are bounded away from 0 and bounded from above, respectively, and do not diverge relative to one another.

4. Error of the normal approximation and model selection

In this section, we establish the asymptotic multivariate normality of the maximum likelihood estimator (MLE) for the data-generating parameter vector $\boldsymbol{\theta}^{\star}$ as its dimension grows. Specifically, we derive a non-asymptotic bound on the quality of the multivariate normal approximation and exhibit scaling conditions on both the model dimension p and the expected number of activated dyads $\mathbb{E} \| \boldsymbol{Y} \|_1$ —under which the approximation error vanishes as the network size tends to infinity. Based on this result, we present a model selection method using multiple hypothesis testing procedures that control the false discovery rate. The main result is presented in Theorem 3, the proof of which is based on a Taylor expansion of the log-likelihood function and through the application of a Lyapunov type bound presented in Raič [33].

In the following, Z will denote a standard multivariate normal random vector, i.e., with mean vector equal to the zero vector and covariance matrix equal to the identity matrix (each of appropriate dimension), and Φ will denote the corresponding probability measure.

Theorem 3. Consider a separable multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ nodes. There exists $N_0 \ge 3$ such that, for all $N \ge N_0$ and any measurable convex set $\mathcal{A} \subseteq \mathbb{R}^p$, the error of the multivariate normal approximation

$$\left|\mathbb{P}((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} \left(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\right) - \Delta \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})\right|$$

is bounded above by

$$\frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}}\sqrt{\frac{p^{7/2}}{\mathbb{E}\,\|\boldsymbol{Y}\|_1}} + \frac{4}{\mathbb{E}\,\|\boldsymbol{Y}\|_1} + \frac{8\,[D_g]^+}{\left(\mathbb{E}\,\|\boldsymbol{Y}\|_1\right)^2}$$

and Δ satisfies

$$\mathbb{P}\left(\|\Delta\|_{2} \leq \frac{\sqrt{2} C^{2} p^{5/2}}{\sqrt{\mathbb{E}\|\boldsymbol{Y}\|_{1}}} \frac{\widetilde{\lambda}_{\max}^{\star}}{(\widetilde{\lambda}_{\min}^{\epsilon})^{5/2}}\right) \geq 1 - \exp\left(-2 p\right) - \frac{5 + 8 C_{0}}{\mathbb{E}\|\boldsymbol{Y}\|_{1}},$$

where C > 0 is the constant given in Theorem 1 and $C_0 > 0$ is the constant given in Assumption 1, both independent of N and p.

Theorem 3 serves as a foundation for establishing the asymptotic normality of the maximum likelihood estimator $\hat{\theta}$. If

$$\lim_{N \to \infty} \left[\frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}} \sqrt{\frac{p^{7/2}}{\mathbb{E} \, \|\boldsymbol{Y}\|_1}} + \frac{4}{\mathbb{E} \, \|\boldsymbol{Y}\|_1} + \frac{8 \left[D_g\right]^+}{\left(\mathbb{E} \, \|\boldsymbol{Y}\|_1\right)^2} \right] = 0,$$

Theorem 3 implies $(I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_1)^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \Delta$ will converge in distribution to a standard multivariate normal random vector, as the error bound on the

multivariate normal approximation will vanish in this case. The term Δ can be viewed as an error term, resulting from the fact that the normal approximation in Theorem 3 is obtained via a multivariate Taylor approximation in order to bridge the distributional gap between key statistics which admit forms amenable to existing theorems for the normal approximation and the parameter vectors of interest, thus introducing an additional source of error in the normal approximation.

While involved, the above condition for asymptotic multivariate normality essentially places restrictions on the dependence induced through the single-layer basis network Y measured by $[D_q]^+$, as well as the smallest eigenvalue of the dyad-based information matrix $I(\boldsymbol{\theta})$ in a neighborhood of the data-generating parameter vector $\boldsymbol{\theta}^{\star}$ as measured by $\widetilde{\lambda}_{\min}^{\epsilon}$, and the model dimension p. As a result, if the information matrix $I(\boldsymbol{\theta})$ is nearly singular at $\boldsymbol{\theta}^{\star}$, in which case $\lambda_{\min}^{\epsilon}$ will be small, the error of the normal approximation will be uniformly larger (all else equal). Likewise, if the edge dependence in \boldsymbol{Y} is large as measured by $[D_a]^+$, we may not have sufficient activated dyads to ensure the error bound is small, as $\|\mathbf{Y}\|_1$ may not be tightly concentrated around $\mathbb{E} \|\mathbf{Y}\|_1$. The dependence of the error approximation on the dimension of the random vector is a known challenge in establishing multivariate normality [e.g., 33]. All quantities which are not explicit constants can increase or decrease with N, with the rates of these increases or decreases having implications for the rate of convergence in distribution. Theorem 3 demonstrates that the allowable scaling for most of quantities is with respect to the expected number of activated dyads $\mathbb{E} \| \mathbf{Y} \|_{1}$.

We further examine Theorem 3 through an example where \mathbf{Y} is a Bernoulli random graph model, which assumes edge variables are independent Bernoulli random variables with probability $\pi \in (0, 1)$. Under this model, $[D_g]^+ = 0$ owing to the independence of edge variables and $\mathbb{E} \| \mathbf{Y} \|_1 = \pi {N \choose 2}$. Under this scenario, we can show that

$$\left| \mathbb{P}((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \Delta \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A}) \right|$$

is bounded above by

$$\frac{166}{\sqrt{\pi \,(\widetilde{\lambda}_{\min}^{\epsilon})^3}} \,\frac{p^{1.75}}{N} + \frac{16}{\pi N^2},$$

with the additional bound

$$\mathbb{P}\left(\|\Delta\|_2 \le \frac{\sqrt{2} C^2 p^{2.5}}{\sqrt{\pi \binom{N}{2}}} \frac{\widetilde{\lambda}_{\max}^{\star}}{(\widetilde{\lambda}_{\min}^{\epsilon})^{2.5}}\right) \ge 1 - \exp\left(-2 p\right) - \frac{5 + 8 C_0}{\pi \binom{N}{2}}$$

where C > 0 is the constant given in Theorem 1 and $C_0 > 0$ is the constant given in Assumption 1, both independent of N and p. If $\tilde{\lambda}_{\min}^{\epsilon}$ and π are both bounded away from 0, then the error of the normal approximation will convergence to 0 provided $(p^{2.5} \tilde{\lambda}_{\max}^{\star}) / N \to 0$ as $N \to \infty$, which is sufficient to ensure $\|\Delta\|_2$ converges in probability to 0. Under the fully saturated model specification for

15

(1) (H = K), the Binomial theorem shows that $p = 2^K - 1 \leq 2^K$. Hence, the dimension restriction on p in turn implies a restriction on the allowable rate of growth of the number of layers K with N, where a sufficient condition for $(p^{2.5} \tilde{\lambda}_{\max}^*) / N \to 0$ is for $K \leq .5 \log N$. In other words, the number of layers K can grow at most logarithmically with N in the fully saturated model. In cases when the number of interaction terms included in the cross-layer dependence probability model is fixed, K may admit a sublinear scaling with N.

4.1. Model selection via univariate testing with FDR control

We outline a procedure for model selection that controls the false discovery rate, leveraging the results of Theorems 1 and 3. Hotelling's *T*-squared statistic can be used to conduct a global test for $H_0: \boldsymbol{\theta}^* = \boldsymbol{\mu}$ versus $H_1: \boldsymbol{\theta}^* \neq \boldsymbol{\mu}$, where $\boldsymbol{\mu} \in \mathbb{R}^p$ is the value of $\boldsymbol{\theta}$ we want to test. We will mostly be interested in the case when $\boldsymbol{\mu} = \mathbf{0}_p$, i.e., the zero vector of dimension *p*. If the global test is rejected, or is not of interest, we can perform model selection by leveraging the multivariate normal approximation to obtain univariate normal approximation results for the components of $\boldsymbol{\hat{\theta}}$ and proceed to test each component: $H_{i,0}: \boldsymbol{\theta}_i^* = \mu_i$ versus $H_{i,1}: \boldsymbol{\theta}_i^* \neq \mu_i$, for $i = 1, \dots p$ and $\mu_i \in \mathbb{R}$. In general, $\mu_i = 0$ will allow us to test whether the estimated effect $\boldsymbol{\hat{\theta}}_i$ is present in the model (i.e., whether $\boldsymbol{\theta}_i^* \neq 0$). One challenge in this approach lies in the fact that the model selection procedure is sensitive to multiple testing error.

To ensure a more reliable procedure for identifying cross-layer dependence effects in multilayer networks while mitigating the risk of spurious discoveries, we elaborate a model selection algorithm that employs suitable multiple testing adjustments to control the false discovery rate. We provide simulation examples of four different univariate testing procedures including Bonferroni, Benjamini-Hochberg, Hochberg, and Holm procedures in Section 5.2. In simulation studies, all four univariate testing procedures exhibit strong statistical power for detecting non-zero parameters, while controlling the false discovery rate at a preset family-wise significance level.

5. Simulation studies

Directly simulating maximum likelihood estimators for network data with dependent edges is challenging because the normalizing constants are often computationally intractable. Computing the normalizing constant requires enumerating all $2^{\binom{N}{2}}$ possible edge combinations for each layer to maximize the true likelihood function. Additionally, dependencies among network dyads prevent factorization of the likelihood, which further complicates direct maximization. As a result, direct maximization of likelihood functions is generally infeasible in these cases. Two predominant methods of approximating the maximum likelihood estimator θ^* when the likelihood function is computationally intractable have emerged in the literature. Monte Carlo maximum likelihood estimation

(MCMLE) [15], which constructs a simulation-based approximation to the likelihood function in order to approximate the maximum likelihood estimator, is an established method for approximating maximum likelihood estimators in the statistical network analysis literature [21]. While able to provide accurate estimates of maximum likelihood estimators for complex models [e.g., 39, 36], a drawback of MCMLE, and other simulation-based estimation methodology, is the computational burden which can scale with both the complexity of the model and the size of the network [5]. In settings where the computation of the MCMLE is impractical, a computationally efficient alternative is provided via the maximum pseudolikelihood estimator (MPLE) [4], whose application to social network analysis and to statistical network analysis dates back to Strauss and Ikeda [40]. Pseudolikelihood-based estimators have the following computational advantages:

- 1. Algorithms are generally deterministic and do not require simulationbased approximation schemes, which aids in reproducibility of results;
- 2. Algorithms are generally more scalable, relative to alternatives such as MCMLE and other simulation-based approximations, and are able to be parallelized to take advantage of larger multicore computing infrastructures which are becoming increasingly common.

In this simulation, we consider the maximum pseudolikelihood estimator, denoted by $\tilde{\theta}$. We conduct simulation studies to investigate the performance of the maximum pseudolikelihood estimator $\tilde{\theta}$ (MPLE), supplementing the theoretical results established in Sections 3 and 4 for maximum likelihood estimators. As will be discussed later in Section 6, we successfully reproduced the sufficient statistics using the MPLE in the application, suggesting that the MPLE for the multilayer network model solves a score equation similar to that of the MLE. This indicates that the MPLE serves as a close approximation and can be a good proxy for the MLE. In section 5.1, we demonstrate the consistency results of Theorem 1 in settings of different data-generating parameters and increasing model dimensions. We conduct simulation studies of the multivariate normal approximation established by Theorem 3 in Section 5.2. Lastly, we discuss several testing procedures for selecting non-zero effects while controlling the false discovery rate (FDR) at a given family-wise significance level α .

In all simulation studies, we sample concordant multilayer networks (X, Y) from (1) with the maximum order of corss-layer interaction H = 2:

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \prod_{\{i,j\} \subset \mathcal{N}} \exp\left(\sum_{k=1}^{K} \theta_k \, x_{i,j}^{(k)} + \sum_{k(6)$$

Unless otherwise specified, the basis network \boldsymbol{Y} is generated from the Bernoulli random graph model.



Fig 2: The relative ℓ_2 -errors between $\tilde{\theta}$ and θ^* decrease as the number of activated dyads increases. Each box is created by 250 replicates of multilayer networks.

5.1. Consistency

The consistency is demonstrated through the decay of the relative ℓ_2 -errors between $\hat{\theta}$ and the data-generating parameter θ^{\star} as the expected number of activated dyads $\mathbb{E} \| \boldsymbol{Y} \|_1$ increases. We generated M = 250 multilayer networks with N = 300 nodes, using M different data-generating parameters. We created these networks for each of ten evenly spaced numbers of activated dyads increasing from 3000 to 30000, and for four different numbers of layers increasing from K = 3 to 6. The model dimension increases from 6 to 21 as K increases from 3 to 6. For each number of activated dyads, number of layers K, and replicate, we sample a multilayer network X from (1) using the specification in (6) with the data-generating parameter vector θ^{\star} populated by randomly selecting each component from the uniform distribution on (-1, 1). We make the exception that components θ_3^{\star} and $\theta_{1,3}^{\star}$ are set to 0. In each replicate, we compute the maximum pseudolikelihood estimator. The results of this simulation study are given in Figure 2, which shows the decay of the relative ℓ_2 -errors between θ and $heta^{\star}$ as the number of activated dyads increases in networks with different number of layers. The broad selection of data-generating parameter values on networks with increasing number of layers verifies that Theorem 1 holds in many practical settings with increasing model dimensions.

5.2. Multivariate normality and model selection

As stated in Section 4 and Theorem 3, the distribution of the maximum likelihood estimator $\hat{\theta}$ converges in distribution to a multivariate normal distribution asymptotically. In order to study the quality of the normal approximation—especially for univariate testing which would be used for the false discovery rate



Fig 3: Q-Q plots and *p*-values of six components of $\tilde{\theta}$ estimated from 250 multilayer network samples at size 1000 on the dense Bernoulli basis network.

TABLE 1 False discovery rates of four procedures for detecting non-zero effects of 6 data-generating parameters $(\theta_1^*, \theta_2^*, \theta_3^*, \theta_4^*, \theta_5^*, \theta_6^*)$ estimated from 250 multilayer network samples at size 1000 on the dense Bernoulli basis network. All FDRs are smaller than .05.

Procedure	$oldsymbol{ heta}_1^\star$	$ heta_2^{\star}$	$ heta_3^{\star}$	$ heta_4^\star$	$ heta_5^{\star}$	$ heta_6^{\star}$
Bonferroni	.004	.002	.001	.002	.001	.005
Benjamini-Hochberg	.014	.014	.014	.011	.017	.020
Hochberg	.012	.008	.009	.008	.011	.016
Holm	.010	.008	.006	.008	.007	.013

control and model selection—we randomly select 6 of the 250 data-generating parameter vectors θ^* used to study the consistency results of Theorem 1 in the simulation study conducted in Section 5.1. We then generate 250 replicates of multilayer network samples by each of these 6 parameter vectors, using specification (6) on four basis network structures with the number of layers K = 3: the Bernoulli random graph model (dense and sparse), the stochastic block model, and the latent space model.

The multivariate normality of $\hat{\theta}$ passed Zhou-Shao's multivariate normal test [46], with *p*-values provided in the Appendix H.1 in the supplement to this paper. We visualize the marginal normality of individual component in $\hat{\theta}$ with a dense Bernoulli basis network in Figure 3, through Q-Q plots of the simulated maximum pseudolikelihood estimators. Univariate tests for normality failed to reject the null hypothesis that each component of $\hat{\theta}$ is marginally normal at a significance level of .05. Additional results studying the multivariate normality of $\hat{\theta}$ on different basis network structures are provided in Appendix H.1 in the supplement to this paper.

We then implement the multiple testing correction procedures of Bonferroni, Benjamini-Hochberg, Hochberg, and Holm, for the 6 selected data-generating

19

ntar g og Bazega o corp		
	Average Node Degree	Number of Edges
Co-Worker Layer	11	378
Advice Layer	5	175
Friendship Layer	5	176

 TABLE 2

 Summary of Lazega's corporate law partnership data with 71 lawyers (nodes).

parameter vectors θ^{\star} with 250 replicates to detect components that are significantly different from 0 while controlling the false discovery rate (FDR) at a family-wise significance level of $\alpha = .05$ —recall that $\theta_{1,3}^{\star}$ and θ_{3}^{\star} of θ^{\star} are set to 0 in each simulation replicate. We estimate the FDR of the four procedures by averaging the false discovery proportions from 250 replicates of each of the 6 randomly selected data-generating parameters θ^{\star} . We provide the estimated FDRs for θ^{\star} on a dense Bernoulli basis network in Table 1. In addition, we show the receiver operating characteristic (ROC) curves for $\tilde{\theta}$ estimating the 6 selected data-generating parameters in each of the subplot of Figure 10, on four basis network structures in Appendix H.2 in the supplement to the paper. Simulation results suggest that the false discovery rate is controlled below the preset threshold α . Different data-generating parameter values affect the tradeoff between the sensitivity and the specificity of the model selection. In general, multilayer networks with a larger effective sample size lead to a larger area under the ROC curve which offers a tool to choose appropriate correction procedures and thresholds for model selection in different scenarios. Additional results on the false discovery rate with different basis network structures are provided in Appendix H.2 in the supplement to the paper.

6. Application

We present a case study using a dataset on corporate law partnership among a Northeastern US corporate law firm in New England collected by Lazega [25]. The dataset collected information about three types of cooperation among 71 lawyers in the corporate law firm, resulting in three networks including the strong-coworker network, the advice network, and the friendship network. Since the cooperation relationship collected are not symmetric, we only consider a connection to be present when both sides acknowledged their cooperation. We treat these three types of networks as a three-layer multilayer network embedded among the 71 lawyers. A summary of this multilayer network is provided in Table 2. We apply the separable multilayer network model in (1) with the specification in (6) to Lazega's lawyer network, i.e., $\boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3, \theta_{1,2}, \theta_{1,3}, \theta_{2,3})$. The maximum pseudolikelihood estimator $\tilde{\boldsymbol{\theta}}$ is computed from the observed network, the results of which are provided in Table 3.

As shown in Table 3, the maximum pseudolikelihood estimates θ_1 , θ_2 , and θ_3 correspond to the estimated single-layer effects of the coworker layer, the advice layer, and the friendship layer, respectively, whereas $\theta_{1,2}$, $\theta_{1,3}$, and $\theta_{2,3}$ correspond to the layer interaction effects. We can calculate the conditional log-



Fig 4: The coworker layer, the advice layer and the friendship layer of Lazega's corporate law partnership network.

 TABLE 3

 MPLEs (and standard errors) of the separable multilayer network model for the Lazega's lawyer network.

$\widetilde{ heta}_1$	$\widetilde{ heta}_2$	$\widetilde{ heta}_3$	$\widetilde{ heta}_{1,2}$	$\widetilde{ heta}_{1,3}$	$\widetilde{ heta}_{2,3}$
-1.450(.263)	-3.334(.244)	-2.695(.256)	1.801(.244)	0.218(.247)	2.458(.231)
Coworker (C)	Advice (A)	Friendship (F)	$C \times A$	$C \times F$	$A \times F$

odds of each edge being present in the multilayer network given the rest of the network. For example, if lawyer i and lawyer j are observed to be coworkers and are friends at the same time, the odds of these two lawyers to have an advice relationship is given by

$$\begin{split} &\frac{\mathbb{P}(X_{i,j}^{(A)} = 1 \mid \boldsymbol{X}_{i,j}^{(C)} = 1, \, \boldsymbol{X}_{i,j}^{(F)} = 1)}{\mathbb{P}(X_{i,j}^{(A)} = 0 \mid \boldsymbol{X}_{i,j}^{(C)} = 1, \, \boldsymbol{X}_{i,j}^{(F)} = 1)} &= \exp\left(\widetilde{\theta}_{2} + \widetilde{\theta}_{1,2} \, x_{i,j}^{(C)} + \widetilde{\theta}_{2,3} \, x_{i,j}^{(F)}\right) \\ &= \exp\left(-3.334 + 1.801 \, x_{i,j}^{(C)} + 2.458 \, x_{i,j}^{(F)}\right) &= 2.522, \end{split}$$

providing interpretation of the interaction and influence among the different layers.

Next, we use the MPLE to reproduce multilayer networks of the same size and compare the sufficient statistics of the simulated networks and the Lazega's lawyer network. We recover the basis network according to Proposition 1, i.e., a dyad is activated if and only if at least one of its layers has a present edge in the Lazega's lawyer multilayer network. We then populate layers of all activated dyads according to equation (6) by the MPLEs obtained in Table 3. Comparisons of the sufficient statistics between the observed Lazega's lawyer network and the simulated networks with 10 replications are provided in Figure 5. Such comparisons serve two key purposes. First, such comparisons are an established method of diagnosing model fit in the statistical network analysis literature [20], and second, provide a check on the approximate solution to the score equation. Note that MPLEs are not guaranteed to reproduce (on average) observed values of sufficient statistics in exponential families—in contrast to MLEs. The relative



Fig 5: Box-plot of reproduced statistics from 10 simulated samples using the MPLE obtained from the Lazega's lawyer network. Red dots are values of the observed sufficient statistics of Lazega's lawyer network.

 ℓ_2 -error of the sufficient statistics between the observed and the average of the 10 simulated networks is 0.09, suggesting a successful re-construction of the observed network statistics.

7. Discussion

In this work, we introduced a flexible class of statistical models for multilayer networks. Key to our approach lies in the integrative nature by which we establish our framework, extending arbitrary strictly positive probability distributions for single-layer networks to multilayer-network models through a novel separable framework with Markov random field specifications. We established the foundations for statistical inference through consistency and multivariate normality results, the results of which have been demonstrated in simulation studies and in an application. The key assumption to our approach lies in the network separability assumption, which necessitates network dyads be conditionally independent given the basis network. This assumption may or may not be valid in practice, which would necessitate the development of generalizations of the framework we established in this work through the relaxation of the conditional independence assumption. Such relaxations would result in more complex dependence structures, requiring new and careful theoretical treatment in order to establish similar statistical foundations of models to the ones we have developed here, representing potential avenues for future research.

Acknowledgements

Jonathan R. Stewart was supported by NSF award SES-2345043 and the Department of Defense Test Resource Management Center under contracts FA807518-D-0002 and FA8075-21-F-0074.

References

- Albert, R. and Barabási, A. L. [2002], 'Statistical mechanics of complex networks', *Reviews of Modern Physics* 74(1), 47.
- [2] Arroyo, J. A., Athreya, A., Cape, J., Chen, G., Priebe, C. E. and Vogelstein, J. T. [2021], 'Inference for multiple heterogeneous networks with a common invariant subspace', *The Journal of Machine Learning Research* 22(1), 6303–6351.
- [3] Athreya, A., Fishkind, D. E., Tang, M., Priebe, C. E., Park, Y., Vogelstein, J. T., Levin, K., Lyzinski, V., Qin, Y. and Sussman, D. L. [2018], 'Statistical inference on random dot product graphs: a survey', *Journal of Machine Learning Research* 18(226), 1–92.
- [4] Besag, J. [1974], 'Spatial interaction and the statistical analysis of lattice systems', Journal of the Royal Statistical Society, Series B 36, 192–225.
- [5] Bhamidi, S., Bresler, G. and Sly, A. [2011], 'Mixing time of exponential random graphs', *The Annals of Applied Probability* 21, 2146–2170.
- [6] Block, P. [2015], 'Reciprocity, transitivity, and the mysterious three-cycle', Social Networks 40, 163–173.
- [7] Brown, L. [1986], Fundamentals of Statistical Exponential Families: With Applications in Statistical Decision Theory, Institute of Mathematical Statistics, Hayworth, CA, USA.
- [8] Butts, C. T. [2020], 'A dynamic process interpretation of the sparse ERGM reference model', *Journal of Mathematical Sociology*.
- [9] Caimo, A. and Gollini, I. [2020], 'A multilayer exponential random graph modelling approach for weighted networks', *Computational Statistics & Data Analysis* 142, 106825.
- [10] Caron, F. and Fox, E. B. [2017], 'Sparse graphs using exchangeable random measures', Journal of the Royal Statistical Society, Series B (with discussion) 79, 1–44.
- [11] Chen, S., Liu, S. and Ma, Z. [2022], 'Global and individualized community detection in inhomogeneous multilayer networks', *The Annals of Statistics* 50(5), 2664–2693.
- [12] Crane, H. and Dempsey, W. [2018], 'Edge exchangeable models for interaction networks', Journal of the American Statistical Association 113(523), 1311–1326.
- [13] Frank, O. [1980], 'Transitivity in stochastic graphs and digraphs', Journal of Mathematical Sociology 7, 199–213.
- [14] Furi, M. and Martelli, M. [1991], 'On the mean value theorem, inequality, and inclusion', *The American Mathematical Monthly* 98(9), 840–846.
- [15] Geyer, C. J. and Thompson, E. A. [1992], 'Constrained Monte Carlo maximum likelihood for dependent data', *Journal of the Royal Statistical Soci*ety, Series B 54, 657–699.
- [16] Hoff, P. D., Raftery, A. E. and Handcock, M. S. [2002], 'Latent space

approaches to social network analysis', Journal of the American Statistical Association **97**, 1090–1098.

- [17] Holland, P. W., Laskey, K. B. and Leinhardt, S. [1983], 'Stochastic block models: some first steps', *Social Networks* 5, 109–137.
- [18] Holland, P. W. and Leinhardt, S. [1972], 'Some evidence on the transitivity of positive interpersonal sentiment', *American Journal of Sociology* 77, 1205–1209.
- [19] Huang, S., Weng, H. and Feng, Y. [2022], 'Spectral clustering via adaptive layer aggregation for multi-layer networks', *Journal of Computational and Graphical Statistics* pp. 1–15.
- [20] Hunter, D. R., Goodreau, S. M. and Handcock, M. S. [2008], 'Goodness of fit of social network models', *Journal of the American Statistical Association* 103, 248–258.
- [21] Hunter, D. R. and Handcock, M. S. [2006], 'Inference in curved exponential family models for networks', *Journal of Computational and Graphical Statistics* 15, 565–583.
- [22] Krivitsky, P. N., Handcock, M. S. and Morris, M. [2011], 'Adjusting for network size and composition effects in exponential-family random graph models', *Statistical Methodology* 8, 319–339.
- [23] Krivitsky, P. N., Handcock, M. S., Raftery, A. E. and Hoff, P. D. [2009], 'Representing degree distributions, clustering, and homophily in social networks with latent cluster random effects models', *Social Networks* **31**, 204– 213.
- [24] Krivitsky, P. N., Koehly, L. M. and Marcum, C. S. [2020], 'Exponentialfamily random graph models for multi-layer networks', *Psychometrika* 85(3), 630–659.
- [25] Lazega, E. [2001], The collegial phenomenon: The social mechanisms of cooperation among peers in a corporate law partnership, Oxford University Press.
- [26] Lei, J., Chen, K. and Lynch, B. [2020], 'Consistent community detection in multi-layer network data', *Biometrika* 107(1), 61–73.
- [27] Li, W., Xu, Y., Yang, J. and Tang, Z. [2012], Finding structural patterns in complex networks, in '2012 IEEE Fifth International Conference on Advanced Computational Intelligence', pp. 23–27.
- [28] Lusher, D., Koskinen, J. and Robins, G. [2013], Exponential Random Graph Models for Social Networks, Cambridge University Press, Cambridge, UK.
- [29] Maathuis, M., Drton, M., Lauritzen, S. and Wainwright, M. [2018], Handbook of graphical models, CRC Press.
- [30] MacDonald, P., Levina, E. and Zhu, J. [2022], 'Latent space models for multiplex networks with shared structure', *Biometrika* 109(3), 683–706.
- [31] McPherson, M., Smith-Lovin, L. and Cook, J. M. [2001], 'Birds of a feather: Homophily in social networks', Annual Review of Sociology 27, 415–444.
- [32] Ortega, J. M. and Rheinboldt, W. C. [2000], Iterative solution of nonlinear equations in several variables, SIAM.
- [33] Raič, M. [2019], 'A multivariate Berry-Esseen theorem with explicit constants', *Bernoulli* 25(4A), 2824–2853.

J. Li et al./Learning cross-layer dependence structure in multilayer networks

24

- [34] Ravikumar, P., Wainwright, M. J. and Lafferty, J. [2010], 'High-dimensional Ising model selection using l₁-regularized logistic regression', *The Annals* of Statistics 38, 1287–1319.
- [35] S. Chen, D. W. a. A. S. [2015], 'Selection and estimation for mixed graphical models', *Biometrika* 102, 47–64.
- [36] Schweinberger, M., Krivitsky, P. N., Butts, C. T. and Stewart, J. [2020], 'Exponential-family models of random graphs: Inference in finite, super, and infinite population scenarios', *Statistical Science* 35, 627–662.
- [37] Sosa, J. and Betancourt, B. [2022], 'A latent space model for multilayer network data', *Computational Statistics & Data Analysis* 169, 107432.
- [38] Stewart, J. R. and Schweinberger, M. [2021], 'Pseudo-likelihood-based *M*estimation of random graphs with dependent edges and parameter vectors of increasing dimension', *arXiv preprint arXiv:2012.07167*.
- [39] Stewart, J., Schweinberger, M., Bojanowski, M. and Morris, M. [2019], 'Multilevel network data facilitate statistical inference for curved ERGMs with geometrically weighted terms', *Social Networks* 59, 98–119.
- [40] Strauss, D. and Ikeda, M. [1990], 'Pseudolikelihood estimation for social networks', Journal of the American Statistical Association 85, 204–212.
- [41] Sundberg, R. [2019], Statistical modelling by exponential families, Vol. 12, Cambridge University Press.
- [42] Vershynin, R. [2018], High-dimensional probability: An introduction with applications in data science, Cambridge University Press, Cambridge, UK.
- [43] Wainwright, M. J. [2019], High-dimensional statistics: A non-asymptotic viewpoint, Vol. 48, Cambridge University Press.
- [44] Wainwright, M. J. and Jordan, M. I. [2008], 'Graphical models, exponential families, and variational inference', Foundations and Trends in Machine Learning 1, 1–305.
- [45] Zhao, P. and Yu, B. [2006], 'On model selection consistency of lasso', Journal of Machine Learning Research.
- [46] Zhou, M. and Shao, Y. [2014], 'A powerful test for multivariate normality', Journal of Applied Statistics 41(2), 351–363.

Supplement: Learning cross-layer dependence structure in multilayer networks

By Jiaheng Li and Jonathan R. Stewart

Department of Statistics, Florida State University

Appendix A: Proof of Proposition 1	1
Appendix B: Proof of Lemma 1	2
Appendix C: Concentration inequalities for multilayer networks	3
Appendix D: Proof of Theorem 1 and Corollary 1	0
Appendix E: Proof of Theorem 2 and Corollary 2 10	0
Appendix F: Proposition 2 and proof	9
Appendix G: Proof of Theorem 3 22	2
Appendix H: Additional simulation results	0

Appendix A: Proof of Proposition 1

We prove Proposition 1 from Section 2. For the first and second results, define the set

 $\mathcal{A}_+ \quad \coloneqq \quad \left\{ (\boldsymbol{x}, \boldsymbol{y}) \in \mathbb{X} \times \mathbb{Y} : h(\boldsymbol{x}, \boldsymbol{y}) = 1 \right\},$

and the vector-valued map $\varphi : \mathbb{X} \mapsto \mathbb{Y}$ by defining its components to be

 $\varphi_{i,j}(\boldsymbol{x}) = \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0), \quad \{i, j\} \subset \mathcal{N},$

populating the vector $\varphi(\boldsymbol{x})$ in the lexicographical ordering of the dyad indices $\{i, j\} \subset \mathbb{N}$. By the definition of $h : \mathbb{X} \times \mathbb{Y} \mapsto \{0, 1\}$ and $\varphi : \mathbb{X} \mapsto \mathbb{Y}, \varphi(\boldsymbol{x}) = \boldsymbol{y}$ for each pair $(\boldsymbol{x}, \boldsymbol{y}) \in \mathcal{A}_+$. Furthermore, the element \boldsymbol{y} is unique for a given $\boldsymbol{x} \in \mathbb{X}$, because if there would exists some $\boldsymbol{y}' \in \mathbb{Y}$ such that $\boldsymbol{y} \neq \boldsymbol{y}'$ with the property that $\{(\boldsymbol{x}, \boldsymbol{y}), (\boldsymbol{x}, \boldsymbol{y}')\} \subseteq \mathcal{A}_+$, then there would exist a pair $\{i, j\} \subset \mathbb{N}$ such that $y_{i,j} = 1 - y'_{i,j}$ which implies $\mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0) \neq y'_{i,j}$. In this case, $h(\boldsymbol{x}, \boldsymbol{y}') = 0$ is contradicting the assumption that $\{(\boldsymbol{x}, \boldsymbol{y})\} \in \mathcal{A}_+$. By (1), the functions f and g are assumed to be strictly positive in their respective domains. Hence, $(\mathbb{X} \times \mathbb{Y}) \setminus \mathcal{A}_+$ is the largest null set of $\mathbb{X} \times \mathbb{Y}$, i.e., $\mathbb{P}_{\theta}(\mathcal{A}) = 0$ if and only if $\mathcal{A} \subseteq (\mathbb{X} \times \mathbb{Y}) \setminus \mathcal{A}_+$. Thus, the first and second results are established.

For the third result, note that g is assumed to be strictly positive on its domain \mathbb{Y} . Hence, $g(\boldsymbol{y}) = \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{Y} = \boldsymbol{y}) > 0$ for all $\boldsymbol{y} \in \mathbb{Y}$ and $\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y})$ is therefore well-defined. By definition,

$$\mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{x}\,|\,oldsymbol{Y}=oldsymbol{y}) \ = \ rac{\mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{x},\,oldsymbol{Y}=oldsymbol{y})}{\mathbb{P}_{oldsymbol{ heta}}(oldsymbol{Y}=oldsymbol{y})}$$

where $\mathbb{P}_{\theta}(Y = y)$ is the marginal probability of event Y = y and is assumed to be equal to g(y). The model form for \mathbb{P}_{θ} given in (1) implies

$$\frac{\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{Y} = \boldsymbol{y})} = \frac{f(\boldsymbol{x}, \boldsymbol{\theta}) g(\boldsymbol{y}) \psi(\boldsymbol{\theta}, \boldsymbol{y})}{g(\boldsymbol{y})} \\ = \exp(\log f(\boldsymbol{x}, \boldsymbol{\theta}) + \log \psi(\boldsymbol{\theta}, \boldsymbol{y})),$$

under the assumption that $h(\boldsymbol{x}, \boldsymbol{y}) = 1$. Hence,

$$\mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{x},oldsymbol{Y}=oldsymbol{y}) \ = \ \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{X}\midoldsymbol{Y}=oldsymbol{y}) \ \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{Y}=oldsymbol{y})$$

so that

$$\log \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{x},oldsymbol{Y}=oldsymbol{y}) \ = \ \log \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}=oldsymbol{X}\midoldsymbol{Y}=oldsymbol{y}) + \log \,g(oldsymbol{y}),$$

as $g(\boldsymbol{y})$ is the marginal probability mass function of \boldsymbol{Y} , i.e., $\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{Y} = \boldsymbol{y}) = g(\boldsymbol{y})$. Lemma 3 establishes that $\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{X} \mid \boldsymbol{Y} = \boldsymbol{y})$ belongs to a minimal exponential family, completing the proof of the third and last result of the proposition.

Appendix B: Proof of Lemma 1

We prove Lemma 1 from Section 2. Using (1),

$$\begin{aligned} -\mathbb{E} \nabla_{\boldsymbol{\theta}}^{2} \ell(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{Y}) &= \sum_{\boldsymbol{y} \in \mathbb{Y}} \sum_{\boldsymbol{x} \in \mathbb{X}} -\nabla_{\boldsymbol{\theta}}^{2} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \,|\, \boldsymbol{Y} = \boldsymbol{y}) \, g(\boldsymbol{y}) \\ &= \sum_{\boldsymbol{y} \in \mathbb{Y}} g(\boldsymbol{y}) \sum_{\boldsymbol{x} \in \mathbb{X}} -\nabla_{\boldsymbol{\theta}}^{2} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \,|\, \boldsymbol{Y} = \boldsymbol{y}) \\ &= \sum_{\boldsymbol{y} \in \mathbb{Y}} g(\boldsymbol{y}) \sum_{\{i, j\} \subset \mathcal{N} : \, y_{i, j} = 1} I(\boldsymbol{\theta}) \\ &= I(\boldsymbol{\theta}) \sum_{\boldsymbol{y} \in \mathbb{Y}} g(\boldsymbol{y}) \,\|\boldsymbol{y}\|_{1} \\ &= I(\boldsymbol{\theta}) \mathbb{E} \|\boldsymbol{Y}\|_{1}. \end{aligned}$$

The above follows by exploiting the conditional independence of vectors $\boldsymbol{x}_{i,j}$ $(\{i, j\} \subset \mathcal{N})$ given $\boldsymbol{Y} = \boldsymbol{y}$ under (1), which implies

$$\ell(oldsymbol{ heta};oldsymbol{x},oldsymbol{y}) = \sum_{\{i,j\} \subset \mathcal{N}} \log \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X}_{i,j}=oldsymbol{x}_{i,j} \,|\, oldsymbol{Y}=oldsymbol{y}),$$

,

and from the fact that the conditional probability distribution of $X_{i,j}$ given Y is a degenerate point mass at **0** when $Y_{i,j} = 0$ so that $-\nabla_{\boldsymbol{\theta}}^2 \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ is a sum of $\|\boldsymbol{y}\|_1$ matrices each equal to $I(\boldsymbol{\theta})$, i.e., given $\boldsymbol{y} \in \mathbb{Y}$, we have

$$\sum_{\boldsymbol{x}\in\mathbb{X}} -\nabla_{\boldsymbol{\theta}}^{2} \,\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \,|\, \boldsymbol{Y} = \boldsymbol{y})$$

$$= \sum_{\{i,j\}\subset\mathbb{N}} \mathbb{E}\left[-\nabla_{\boldsymbol{\theta}}^{2} \,L_{i,j}(\boldsymbol{\theta}, \boldsymbol{X}_{i,j}, \boldsymbol{Y}) \,|\, \boldsymbol{Y} = \boldsymbol{y}\right] = \sum_{\{i,j\}\subset\mathbb{N}: \, y_{i,j} = 1} I(\boldsymbol{\theta}).$$

The fact that $I(\boldsymbol{\theta})$ is constant for all pairs $\{i, j\} \subset \mathbb{N}$ satisfying $Y_{i,j} = 1$ follows from the form of (1), which assumes each vector $\boldsymbol{X}_{i,j}$ ($\{i, j\} \subset \mathbb{N}$) is conditionally independent and identically distributed, conditional on \boldsymbol{Y} . Hence,

$$\mathbb{E}\left[-\nabla_{\boldsymbol{\theta}}^{2}\,\ell(\boldsymbol{\theta};\boldsymbol{X},\boldsymbol{Y})\right] = I(\boldsymbol{\theta})\,\mathbb{E}\,\|\boldsymbol{Y}\|_{1},$$

which in turn implies

$$\begin{split} \lambda_{\min}(-\mathbb{E}\nabla_{\boldsymbol{\theta}}^{2}\,\ell(\boldsymbol{\theta};\boldsymbol{X},\boldsymbol{Y})) &= \lambda_{\min}(I(\boldsymbol{\theta}))\,\mathbb{E}\,\|\boldsymbol{Y}\|_{1}\\ \lambda_{\max}(-\mathbb{E}\nabla_{\boldsymbol{\theta}}^{2}\,\ell(\boldsymbol{\theta};\boldsymbol{X},\boldsymbol{Y})) &= \lambda_{\max}(I(\boldsymbol{\theta}))\,\mathbb{E}\,\|\boldsymbol{Y}\|_{1}. \end{split}$$

г		
L		
L		

Appendix C: Concentration inequalities for multilayer networks

We establish the concentration inequality of gradients of log-likelihood functions of multilayer networks in Lemma 2. Recall the definition $[D_g]^+ := \max\{0, D_g\}$, where

$$D_g := \sum_{\{i,j\} \prec \{v,w\} \subset \mathcal{N}} \mathbb{C}(Y_{i,j}, Y_{v,w}),$$

with $\{i, j\} \prec \{v, w\}$ implying the sum is taken with respect to the lexicographical ordering of pairs of nodes, and where $g : \mathbb{Y} \mapsto (0, 1)$ is the marginal probability mass function of Y.

Lemma 2. Consider a multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ nodes and $K \ge 1$ layers. Define $\nabla_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{y}) \coloneqq -\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$, where $\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ is the log-likelihood function. Then, for all t > 0and $\boldsymbol{\theta} \in \mathbb{R}^p$, the probability

$$\mathbb{P}\left(\|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) - \mathbb{E}\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y})\|_{2} \ge t\right)$$

is bounded above by

$$\exp\left(-\frac{t^2}{36\,\widetilde{\lambda}^{\star}_{\max}\left(\mathbb{E}\,\|\boldsymbol{Y}\|_1 + [D_g]^+\right) + 2\,\sqrt{p}\,t} + \log\,p\right) + \frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_1}$$

PROOF OF LEMMA 2. By Proposition 1,

$$\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) = \log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) + \log g(\boldsymbol{y})$$

Thus,

$$-\nabla_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{y}) = \nabla_{\boldsymbol{\theta}} \log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) + \nabla_{\boldsymbol{\theta}} \log g(\boldsymbol{y})$$

= $s(\boldsymbol{x}) - \mathbb{E}_{\boldsymbol{\theta}} s(\boldsymbol{X}),$ (7)

as $g(\boldsymbol{y}) = \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{Y} = \boldsymbol{y})$ is assumed to not be a function of $\boldsymbol{\theta}$. The last equation in (7) follows from Lemma 3, which showed that $\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y})$ is a minimal exponential family with the natural parameter vector $\boldsymbol{\theta} \in \mathbb{R}^p$ and the sufficient statistic vector $\boldsymbol{s}(\boldsymbol{x})$ defined in Lemma 3, inserting the familiar form of the score equation of an exponential family with respect to the natural parameter vector [e.g., Proposition 3.10, p. 32, 41]. Thus,

$$egin{aligned} -(
abla_{m{ heta}}(m{X},m{Y})-\mathbb{E}\,
abla_{m{ heta}}(m{X},m{Y})) &=& s(m{X})-\mathbb{E}_{m{ heta}}\,s(m{X})-\mathbb{E}\,[s(m{X})-\mathbb{E}_{m{ heta}}\,s(m{X})] \ &=& s(m{X})-\mathbb{E}\,s(m{X}). \end{aligned}$$

Let t > 0 and $\boldsymbol{\theta} \in \mathbb{R}^p$ be arbitrary and fixed and define $\mathcal{D}_2(\boldsymbol{\theta}, t)$ to be the event that $\|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) - \mathbb{E} \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y})\|_2 \ge t$, i.e.,

$$\mathcal{D}_2(\boldsymbol{ heta},t) = \{ \boldsymbol{x} \in \mathbb{X} : \| \boldsymbol{s}(\boldsymbol{x}) - \mathbb{E} \, \boldsymbol{s}(\boldsymbol{X}) \|_2 \geq t \}.$$

Let $\epsilon > 0$ and define $\mathcal{E}(\epsilon)$ to be the event that $|||\mathbf{Y}||_1 - \mathbb{E}||\mathbf{Y}||_1| \le \epsilon$, i.e.,

$$\mathcal{E}(\epsilon) = \{ \boldsymbol{y} \in \mathbb{Y} : |||\boldsymbol{y}||_1 - \mathbb{E}||\boldsymbol{Y}||_1| \le \epsilon \}.$$

We assume that $\epsilon > 0$ is chosen so that $\mathcal{E}(\epsilon)$ is not empty, which implies $\mathbb{P}(\mathcal{E}(\epsilon)) > 0$ as $g(\boldsymbol{y})$ is assumed to be strictly positive on \mathbb{Y} . By the law of total probability,

$$\mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta}, t)) = \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta}, t) | \mathcal{E}(\epsilon)) \mathbb{P}(\mathcal{E}(\epsilon)) + \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta}, t) | \mathcal{E}(\epsilon)^{c}) \mathbb{P}(\mathcal{E}(\epsilon)^{c}) \\
\leq \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta}, t) | \mathcal{E}(\epsilon)) + \mathbb{P}(\mathcal{E}(\epsilon)^{c}).$$
(8)

Note that we have not necessarily guaranteed that $\mathbb{P}(\mathcal{E}(\epsilon)^c) > 0$. However, if $\mathbb{P}(\mathcal{E}(\epsilon)^c) = 0$ the non-conditional form of the law of total probability would yield the bound

$$\mathbb{P}\left(\mathcal{D}_2(\boldsymbol{\theta}, t)\right) \leq \mathbb{P}\left(\mathcal{D}_2(\boldsymbol{\theta}, t) \,|\, \mathcal{E}(\epsilon)\right),$$

which is strictly sharper than the bound we give in (8). We will use a divideand-conquer strategy to bound each probability in (8) in turn.

To bound the first term in (8), let $\mathcal{U} := \{ \boldsymbol{u} \in \mathbb{R}^p : \|\boldsymbol{u}\|_2 \leq 1 \}$ be a closed unit ball in \mathbb{R}^p . Define an ϵ -net \mathcal{V}_{ϵ} of $\mathcal{U} \subset \mathbb{R}^p$. By Corollary 4.2.13 of [42], there exists an ϵ -net $\mathcal{V}_{\epsilon} \subset \mathcal{U}$ such that its cardinality satisfies log $|\mathcal{V}_{\epsilon}| \leq p \log (2\epsilon^{-1} + 1)$.

4

Taking $\epsilon = 1/2$, for each $\boldsymbol{u} \in \mathcal{U}$, there exists a $\boldsymbol{v} \in \mathcal{V}_{1/2}$ such that $\|\boldsymbol{u} - \boldsymbol{v}\|_2 \leq 1/2$, and by the Cauchy-Schwarz inequality,

$$\begin{array}{ll} \langle \boldsymbol{u}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\rangle &= \langle \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\rangle + \langle \boldsymbol{u} - \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\rangle \\ &\leq \langle \boldsymbol{u}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\rangle + \| \boldsymbol{u} - \boldsymbol{v} \|_2 \, \| \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \|_2 \quad (9) \\ &\leq \langle \boldsymbol{u}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \,\rangle + \frac{1}{2} \, \| \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \|_2. \end{array}$$

If $\|\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\|_2 \neq 0$, we can choose

$$u_i = \frac{\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})_i}{\|\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\|_2},$$

so that $\|\boldsymbol{u}\|_2 \leq 1$ and $\boldsymbol{u} \in \mathcal{U}$. Next, re-write

$$\begin{split} \langle \boldsymbol{u} , \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \, \rangle &= \frac{1}{\| \langle \boldsymbol{u} , \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \, \rangle \|_2} \, \sum_{i=1}^p \, (\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})_i)^2 \\ &= \| \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \|_2, \end{split}$$

and together with (9), we have

$$\|\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\|_{2} \leq 2 \max_{\boldsymbol{v} \in \mathcal{V}_{1/2}} \langle \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \rangle.$$
(10)

If $\|\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\|_2 = 0$, the inequality (10) holds trivially. As a result of (10), for any t > 0,

$$\begin{split} \mathbb{P}\left(\mathfrak{D}_{2}(\boldsymbol{\theta},t) \left| \left. \mathcal{E}(\epsilon) \right) \right| &\leq \mathbb{P}\left(2 \max_{\boldsymbol{v} \in \mathcal{V}_{1/2}} \left\langle \boldsymbol{v} \,, \, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta};\boldsymbol{x},\boldsymbol{y}) \right\rangle \geq t \right) \\ &\leq \sum_{\boldsymbol{v} \in \mathcal{V}_{1/2}} \mathbb{P}\left(\left\langle \boldsymbol{v} \,, \, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta};\boldsymbol{x},\boldsymbol{y}) \right\rangle \geq \frac{t}{2} \right) \\ &\leq \exp\left(p \log 5 \right) \max_{\boldsymbol{v} \in \mathcal{V}_{1/2}} \mathbb{P}\left(\left\langle \boldsymbol{v} \,, \, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta};\boldsymbol{x},\boldsymbol{y}) \right\rangle \geq \frac{t}{2} \right). \end{split}$$

The last inequality is true because log $|\mathcal{V}_{1/2}| \leq p$ log 5. Note that

$$\langle \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \rangle = \sum_{l=1}^{p} v_{l} [\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})]_{l}$$

$$= \sum_{l=1}^{p} v_{l} [s_{l}(\boldsymbol{x}) - \mathbb{E} s_{l}(\boldsymbol{X})].$$

$$(11)$$

The form of (1) implies, through factorization principles, that the dyad-based vectors $\mathbf{X}_{i,j}$ ($\{i, j\} \subset \mathbb{N}$) are conditionally independent given \mathbf{Y} [e.g., 29, p.

11–13]. Hence, using Lemma 3, the components of the sufficient statistic vector decompose into the sum

$$s_l(\boldsymbol{X}) = \sum_{\{i,j\} \subset \mathcal{N}} s_{l,i,j}(\boldsymbol{X}_{i,j}), \quad l \in \{1,\ldots,p\},$$

so that the components of s(X) are sums of bounded conditionally independent random variables given Y. As a result, equation (11) can be further decomposed into sums of independent random variables:

$$\langle \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \rangle = \sum_{\{i,j\} \subset \mathcal{N}} \sum_{l=1}^{p} v_l [s_{l,i,j}(\boldsymbol{x}_{i,j}) - \mathbb{E} s_{l,i,j}(\boldsymbol{X}_{i,j})].$$

Using the forms for $s_{l,i,j}(\mathbf{X}_{i,j})$ given in Lemma 3, we have $0 \leq s_{l,i,j}(\mathbf{X}_{i,j}) \leq Y_{i,j}$ \mathbb{P} -almost surely, because $s_{l,i,j}(\mathbf{X}_{i,j}) \in \{0,1\}$ and $s_{l,i,j}(\mathbf{X}_{i,j}) = 0$ if $Y_{i,j} = 0$ \mathbb{P} almost surely. Then for each $\{i, j\} \subset \mathbb{N}$, we have

$$\mathbb{E}\sum_{l=1}^{p} v_l \left[s_{l,i,j}(\boldsymbol{x}_{i,j}) - \mathbb{E} s_{l,i,j}(\boldsymbol{X}_{i,j}) \right] = 0,$$

and by the Cauchy-Schwarz inequality, we obtain

$$\left| \sum_{l=1}^{p} v_{l} \left[s_{l,i,j}(\boldsymbol{x}_{i,j}) - \mathbb{E} \, s_{l,i,j}(\boldsymbol{X}_{i,j}) \right] \right| \leq \| \boldsymbol{v} \|_{2} \sqrt{p} \, \| \boldsymbol{s}_{l,i,j}(\boldsymbol{x}_{i,j}) - \mathbb{E} \, \boldsymbol{s}_{l,i,j}(\boldsymbol{X}_{i,j}) \|_{\infty} \\ \leq \frac{3}{2} \sqrt{p}.$$

The last inequality follows from

$$\|m{v}\|_2 \leq \|m{u}\|_2 + \|m{u} - m{v}\|_2 \leq 1 + rac{1}{2} \leq rac{3}{2}.$$

The inequality is true because the construction of the ϵ -net $\mathcal{V}_{1/2} \subset \mathcal{U}$ with $\epsilon = 1/2$ ensures that such a $\boldsymbol{u} \in \mathcal{U}$ exists. We next bound the variance by

$$\begin{aligned} \mathbb{V} & \sum_{\{i,j\} \subset \mathbb{N}} \sum_{l=1}^{p} v_l \left[s_{l,i,j}(\boldsymbol{x}_{i,j}) - \mathbb{E} s_{l,i,j}(\boldsymbol{X}_{i,j}) \right] \\ &= \sum_{\{i,j\} \subset \mathbb{N}} \sum_{m=1}^{p} \sum_{n=1}^{p} \mathbb{C} \left(v_m \, s_{m,i,j}(\boldsymbol{X}_{i,j}) \,, \, v_n \, s_{n,i,j}(\boldsymbol{X}_{i,j}) \right) \\ &= \sum_{\{i,j\} \subset \mathbb{N}} \sum_{m=1}^{p} \sum_{n=1}^{p} v_m \, v_n \, \mathbb{C} (s_{m,i,j}(\boldsymbol{X}_{i,j}) \, s_{n,i,j}(\boldsymbol{X}_{i,j})) \\ &= \langle \boldsymbol{v} \,, \, \| \boldsymbol{y} \|_1 \, I(\boldsymbol{\theta}^*) \, \boldsymbol{v} \rangle \\ &\leq \| \boldsymbol{v} \|_2^2 \, \| \boldsymbol{y} \|_1 \, \widetilde{\lambda}^*_{\max} \\ &\leq \frac{9}{4} \, \| \boldsymbol{y} \|_1 \, \widetilde{\lambda}^*_{\max}, \end{aligned}$$

where $\tilde{\lambda}_{\max}^{\star}$ is the largest eigenvalue of the Fisher information of individual activated dyad defined in Lemma 1 evaluated at the data-generating parameter θ^{\star} . We then apply the one-sided Bernstein's inequality to obtain the upper bound for the conditional probability of $\mathcal{D}_2(\theta, t)$ as follows: [e.g., 42, Theorem 2.8.4]

$$\mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) | \boldsymbol{Y} = \boldsymbol{y}) \leq \exp(p \log 5) \max_{\boldsymbol{v} \in \mathcal{V}_{1/2}} \mathbb{P}\left(\langle \boldsymbol{v}, \nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \rangle \geq \frac{t}{2}\right) \\
\leq \exp\left(\frac{-\frac{(t/2)^{2}}{2}}{\frac{9}{4} \|\boldsymbol{y}\|_{1} \, \tilde{\lambda}_{\max}^{\star} + \frac{1}{3} \, \frac{3}{2} \sqrt{p} \frac{t}{2}} + p \log 5\right) \qquad (12) \\
= \exp\left(\frac{-t^{2}}{18 \|\boldsymbol{y}\|_{1} \, \tilde{\lambda}_{\max}^{\star} + 2 \sqrt{p} t} + p \log 5\right).$$

Using the law of total probability, we bound $\mathbb{P}(\mathcal{D}_2(\boldsymbol{\theta}, t) | \mathcal{E}(\epsilon))$ as follows:

$$\mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) | \mathcal{E}(\epsilon)) = \sum_{\boldsymbol{y} \in \mathbb{Y}} \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) \cap [\boldsymbol{Y} = \boldsymbol{y}] | \mathcal{E}(\epsilon))$$

$$= \sum_{\boldsymbol{y} \in \mathcal{E}(\epsilon)} \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) \cap [\boldsymbol{Y} = \boldsymbol{y}] | \mathcal{E}(\epsilon))$$

$$= \sum_{\boldsymbol{y} \in \mathcal{E}(\epsilon)} \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) | [\boldsymbol{Y} = \boldsymbol{y}] \cap \mathcal{E}(\epsilon)) \mathbb{P}(\boldsymbol{Y} = \boldsymbol{y} | \mathcal{E}(\epsilon))$$

$$= \sum_{\boldsymbol{y} \in \mathcal{E}(\epsilon)} \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) | \boldsymbol{Y} = \boldsymbol{y}) \frac{\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}(\mathcal{E}(\epsilon))},$$
(13)

noting that $[\mathbf{Y} = \mathbf{y}] \cap \mathcal{E}(\epsilon) = [\mathbf{Y} = \mathbf{y}]$ whenever $\mathbf{y} \in \mathcal{E}(\epsilon)$ and in the case when $\mathbf{y} \notin \mathcal{E}(\epsilon)$, the intersection is empty, implying

$$\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y} \,|\, \mathcal{E}(\epsilon)) \quad = \quad \frac{\mathbb{P}([\boldsymbol{Y} = \boldsymbol{y}] \cap \mathcal{E}(\epsilon))}{\mathbb{P}(\mathcal{E}(\epsilon))} \quad = \quad \begin{cases} \frac{\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}(\mathcal{E}(\epsilon))} & \boldsymbol{y} \in \mathcal{E}(\epsilon) \\ 0 & \boldsymbol{y} \notin \mathcal{E}(\epsilon) \end{cases}.$$

We now bound (13) using the bound in (12):

$$\begin{split} &\sum_{\boldsymbol{y}\in\mathcal{E}(\epsilon)} \mathbb{P}(\mathcal{D}_{2}(\boldsymbol{\theta},t) \,|\, \boldsymbol{Y} = \boldsymbol{y}) \, \frac{\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}(\mathcal{E}(\epsilon))} \\ &\leq \sum_{\boldsymbol{y}\in\mathcal{E}(\epsilon)} \exp\left(\frac{-t^{2}}{18 \,\|\boldsymbol{y}\|_{1} \,\widetilde{\lambda}_{\max}^{\star} + 2 \sqrt{p} \,t} \,+ p \,\log 5\right) \, \frac{\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}(\mathcal{E}(\epsilon))} \\ &\leq \exp\left(\frac{-t^{2}}{18 \,(\mathbb{E} \,\|\boldsymbol{Y}\|_{1} \,+ \epsilon) \,\widetilde{\lambda}_{\max}^{\star} + 2 \sqrt{p} \,t} \,+ p \,\log 5\right) \, \sum_{\boldsymbol{y}\in\mathcal{E}(\epsilon)} \, \frac{\mathbb{P}(\boldsymbol{Y} = \boldsymbol{y})}{\mathbb{P}(\mathcal{E}(\epsilon))} \\ &= \exp\left(\frac{-t^{2}}{18 \,(\mathbb{E} \,\|\boldsymbol{Y}\|_{1} \,+ \epsilon) \,\widetilde{\lambda}_{\max}^{\star} \,+ 2 \sqrt{p} \,t} \,+ p \,\log 5\right), \end{split}$$

showing

$$\mathbb{P}\left(\mathcal{D}_{2}(\boldsymbol{\theta},t) \,|\, \mathcal{E}(\epsilon)\right) \leq \exp\left(\frac{-t^{2}}{18\left(\mathbb{E} \,\|\boldsymbol{Y}\|_{1} + \epsilon\right)\widetilde{\lambda}_{\max}^{\star} + 2\sqrt{p}t} + p \log 5\right)$$

The replacement of $\|\boldsymbol{y}\|_1$ by $\mathbb{E}\|\boldsymbol{Y}\|_1 + \epsilon$ follows because $\|\boldsymbol{y}\|_1 \leq \mathbb{E}\|\boldsymbol{Y}\|_1 + \epsilon$ for $\boldsymbol{y} \in \mathcal{E}(\epsilon)$, resulting in the upper bound above. We bound the second term in the inequality (8) using Chebyshev's inequality:

$$\begin{aligned} \mathbb{P}(\mathcal{E}(\epsilon)^c) &= \mathbb{P}(|\|\boldsymbol{Y}\|_1 - \mathbb{E} \, \|\boldsymbol{Y}\|_1| > \epsilon) \\ &\leq \mathbb{P}(|\|\boldsymbol{Y}\|_1 - \mathbb{E} \, \|\boldsymbol{Y}\|_1| \ge \epsilon) \\ &\leq \frac{\mathbb{V}(\|\boldsymbol{Y}\|_1)}{\epsilon^2}. \end{aligned}$$

We bound the variance $\mathbb{V}(\|\boldsymbol{Y}\|_1)$ as follows:

$$\begin{aligned} \mathbb{V}(\|\boldsymbol{Y}\|_{1}) &= \sum_{\{i,j\}\subset\mathcal{N}} \mathbb{V}Y_{i,j} + 2\sum_{\{i,j\}\prec\{v,w\}\subset\mathcal{N}} \mathbb{C}(Y_{i,j}, Y_{v,w}) \\ &\leq \mathbb{E}\|\boldsymbol{Y}\|_{1} + 2\sum_{\{i,j\}\prec\{v,w\}\subset\mathcal{N}} \mathbb{C}(Y_{i,j}, Y_{v,w}), \end{aligned}$$

noting $Y_{i,j} \in \{0,1\}$ so that $\mathbb{V} Y_{i,j} = \mathbb{P}(Y_{i,j} = 1) \mathbb{P}(Y_{i,j} = 0) \leq \mathbb{E} Y_{i,j}$. Hence,

$$\mathbb{P}(\mathcal{E}(\epsilon)^{c}) \leq \frac{\mathbb{E} \|\boldsymbol{Y}\|_{1} + 2\sum_{\{i,j\}\prec\{v,w\}\subset\mathbb{N}} \mathbb{C}(Y_{i,j}, Y_{v,w})}{\epsilon^{2}} \\
= \frac{\mathbb{E} \|\boldsymbol{Y}\|_{1} + 2[D_{g}]^{+}}{\epsilon^{2}}.$$
(14)

Taking $\epsilon = \mathbb{E} \left\| \boldsymbol{Y} \right\|_1 + 2 \left[D_g \right]^+ > 0$ shows that $\mathbb{P}(\mathcal{E}(\epsilon)^c) \leq (\mathbb{E} \left\| \boldsymbol{Y} \right\|_1)^{-1}$ and

$$\mathbb{P}\left(\mathcal{D}_{2}(\boldsymbol{\theta},t) \,|\, \mathcal{E}(\epsilon)\right) \leq \exp\left(\frac{-t^{2}}{36\,\widetilde{\lambda}_{\max}^{\star}\left(\mathbb{E}\,\|\boldsymbol{Y}\|_{1}+[D_{g}]^{+}\right)\,+\,2\,\sqrt{p}\,t}\,+\,p\,\log\,5\right).$$

Combining all results shows that

$$\mathbb{P}\left(\|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) - \mathbb{E}\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y})\|_{2} \ge t\right)$$

is bounded above by

$$\exp\left(\frac{-t^2}{36\,\widetilde{\lambda}_{\max}^{\star}\left(\mathbb{E}\,\|\boldsymbol{Y}\|_1+[D_g]^+\right)\,+\,2\,\sqrt{p}\,t}\,+\,p\,\log\,5\right)\,+\,\frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_1}.$$

As a final matter, note that this choice of $\epsilon > 0$ ensures $\mathcal{E}(\epsilon)$ contains all $\boldsymbol{y} \in \mathbb{Y}$ with $\|\boldsymbol{y}\|_1 \in [0, 2(\mathbb{E} \|\boldsymbol{Y}\| + [D_g]^+)]$ as the empty graph is an element of \mathbb{Y} with 0 edges.

8

C.1. Auxiliary results

Lemma 3. Consider a multilayer network model following the form of equation (1) with maximum interaction term $H \leq K$ and is defined on a set of $N \geq 3$ nodes and $K \geq 1$ layers. Then the following hold:

1. The conditional probability mass function of X given Y is an exponential family:

$$\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) \propto h(\boldsymbol{x}, \boldsymbol{y}) \exp\left(\langle \boldsymbol{\theta}, \boldsymbol{s}(\boldsymbol{x}) \rangle\right),$$

where

$$h(\boldsymbol{x}, \, \boldsymbol{y}) = \prod_{\{i,j\} \subset \mathcal{N}} \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0)^{y_{i,j}} \, \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 = 0)^{1-y_{i,j}},$$

the sufficient statistic vector $s : \mathbb{X} \mapsto \mathbb{R}^p$ and the natural parameter vector $\boldsymbol{\theta} \in \mathbb{R}^p$.

2. For each $l \in \{1, ..., p\}$, there exists $h \in \{1, ..., H\}$ and $\{k_1, ..., k_h\} \subseteq \{1, ..., K\}$ such that the l^{th} component of the sufficient statistic vector $s(\boldsymbol{x})$ can be written as

$$s_l(\boldsymbol{x}) = \sum_{\{i,j\} \subset \mathbb{N}} s_{l,i,j}(\boldsymbol{x}) = \prod_{r=1}^h x_{i,j}^{(k_r)}.$$
 (15)

3. The exponential family outlined above is minimal, full, and regular.

PROOF OF LEMMA 3. First, the form of the conditional probability distribution of X given Y derived in Proposition 1 is given by

$$\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \,|\, \boldsymbol{Y} = \boldsymbol{y}) = \exp\left(\log f(\boldsymbol{x}, \boldsymbol{\theta}) + \log \psi(\boldsymbol{\theta}, \boldsymbol{y})\right), \quad (16)$$

provided $h(\boldsymbol{x}, \boldsymbol{y}) = 1$. The form of (1) suggests that (16) will be a minimal exponential family in canonical form due to the form of the Markov random field specification for $f(\boldsymbol{\theta}, \boldsymbol{x})$ and the definition of $\psi(\boldsymbol{\theta}, \boldsymbol{y})$. From the form of $f(\boldsymbol{x}, \boldsymbol{\theta})$ in (1),

$$\log f(\boldsymbol{x}, \boldsymbol{\theta}) = \sum_{\{i, j\} \subset \mathcal{N}} \left(\sum_{k=1}^{K} \theta_k x_{i, j}^{(k)} + \sum_{k < l}^{K} \theta_{k, l} x_{i, j}^{(k)} x_{i, j}^{(l)} + \ldots + \sum_{k_1 < \ldots < k_H}^{K} \theta_{k_1, k_2, \ldots, k_H} x_{i, j}^{(k_1)} \cdots x_{i, j}^{(k_H)} \right),$$

where $H \leq K$ is the highest order of cross-layer interactions included in the model. We write $\theta_{k_1,k_2,...,k_h}$ to reference the *h*-order interaction parameter for the interaction term among layers $\{k_1,...,k_h\} \subseteq \{1,...,K\}$. As specified, $\psi(\boldsymbol{\theta}, \boldsymbol{y})$ is the normalizing constant for the exponential family. As such, the natural parameter space of the exponential family is \mathbb{R}^p as the support of \mathbb{X} is finite, which implies $\psi(\boldsymbol{\theta}, \boldsymbol{y}) < \infty$ for all $\boldsymbol{\theta} \in \mathbb{R}^p$ and $\boldsymbol{y} \in \mathbb{Y}$. We establish minimality by noting that the components of the parameter vector $\boldsymbol{\theta}$ satisfy no linear or affine

constraints. Attached to each parameter θ_{k_1,\ldots,k_h} ($\{k_1,\ldots,k_h\} \subseteq \{1,\ldots,K\}$, $h \in \{1,\ldots,H\}$) is the sufficient statistic

$$s_{k_1,...,k_h}(\boldsymbol{x}) = \sum_{\{i,j\} \subset \mathcal{N}} x_{i,j}^{(k_1)} \cdots x_{i,j}^{(k_h)}.$$

Each statistic s_{k_1,\ldots,k_h} is a function of distinct, non-degenerate random variables, provided $\|\boldsymbol{y}\|_1 > 0$, and so none of the statistics s_{k_1,\ldots,k_h} satisfy any linear or affine constraints. Hence, (1) specifies a minimal and full exponential family with natural parameter space \mathbb{R}^p of dimension $p = \sum_{h=1}^{H} {K \choose h}$ and sufficient statistic vector $s(\boldsymbol{x})$ with components $s_{k_1,\ldots,k_h}(\boldsymbol{x})$ ($\{k_1,\ldots,k_h\} \subseteq \{1,\ldots,K\}, h =$ $1,\ldots,H$). Regularity follows trivially [e.g., Proposition 3.7, pp. 28, 41]. The form of (15) outlines this for a linear indexing of the components of the sufficient statistic vector.

Appendix D: Proof of Theorem 1

We prove Theorem 1 from Section 3. By Proposition 1, observing X = x implies we observe Y = y, as for each given $x \in \mathbb{X}$, Y = y (P-a.s.) for one and only one $y \in \mathbb{Y}$ given by

$$y_{i,j} = \mathbb{1}(\|\boldsymbol{x}_{i,j}\|_1 > 0), \quad \{i,j\} \subset \mathcal{N}.$$

Denote the gradient of $-\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ by

$$abla_{oldsymbol{ heta}}(oldsymbol{x},oldsymbol{y}) \ \coloneqq \ -
abla_{oldsymbol{ heta}}\,\ell(oldsymbol{ heta};oldsymbol{x},oldsymbol{y})$$

and the expected Hessian matrix of the negative log-likelihood by

$$H(\theta) := -\mathbb{E} \nabla^2_{\theta} \ell(\theta; X, Y).$$

Theorem 6.3.4 of Ortega and Rheinboldt [32] states that if

$$(\boldsymbol{ heta} - \boldsymbol{ heta}^{\star})^{ op} \nabla_{\boldsymbol{ heta}}(\boldsymbol{x}, \boldsymbol{y}) \geq 0$$

for all $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \epsilon)$, where $\partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \epsilon)$ is the boundary of the set

$$\mathcal{B}_2(\boldsymbol{ heta}^\star,\epsilon) = \{ \boldsymbol{ heta} \in \mathbb{R}^p : \| \boldsymbol{ heta} - \boldsymbol{ heta}^\star \|_2 < \epsilon \},$$

then $\nabla_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{y})$ has a root in $\mathcal{B}_2(\boldsymbol{\theta}^*, \epsilon) \cup \partial \mathcal{B}_2(\boldsymbol{\theta}^*, \epsilon)$, i.e., $\widehat{\boldsymbol{\theta}}$ exists and satisfies $\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \leq \epsilon$. Note that a root of $\nabla_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{y})$ is also a root of $-\nabla_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{y})$; in what follows, we consider finding a maximizer of $\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ by finding a minimizer of $-\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$. The classification of roots as maximizers/minimizers is justified from the fact that that $\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ is concave in $\boldsymbol{\theta}$, a fact which follows from Proposition 1, as $g(\boldsymbol{y})$ is constant in $\boldsymbol{\theta}$ and $\log \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} | \boldsymbol{Y} = \boldsymbol{y})$ is the

10

11

log-likelihood of a minimal, full, and regular exponential family with natural parameter vector $\boldsymbol{\theta}$ and thus is strictly concave in $\boldsymbol{\theta}$ [Proposition 3.10, p. 32, 41]. By the multivariate mean-value theorem [14, Theorem 5],

$$\begin{aligned} (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \mathbb{E} \, \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) &= (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \mathbb{E} \nabla_{\boldsymbol{\theta}^{\star}}(\boldsymbol{X}, \boldsymbol{Y}) + (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \boldsymbol{H}(\dot{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}) \\ &= (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \boldsymbol{H}(\dot{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}), \end{aligned}$$

where $\dot{\boldsymbol{\theta}} = t \boldsymbol{\theta} + (1-t) \boldsymbol{\theta}^{\star}$ (some $t \in [0,1]$) and by invoking Lemma 2 of Stewart and Schweinberger [38], which shows that both the expected log-likelihood and log-pseudolikelihood of a minimal exponential family is uniquely maximized at the data-generating parameter vector $\boldsymbol{\theta}^{\star}$, implying $\mathbb{E} \nabla_{\boldsymbol{\theta}^{\star}}(\boldsymbol{X}, \boldsymbol{Y}) = 0$. Let $\gamma \in (0, \epsilon)$ and arbitrarily take $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \gamma)$. Then

$$\begin{aligned} (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \boldsymbol{H}(\dot{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}) &= \frac{(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \boldsymbol{H}(\dot{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})}{(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|_{2}^{2} \\ &\geq \gamma^{2} \lambda_{\min}(\boldsymbol{H}(\dot{\boldsymbol{\theta}})), \end{aligned}$$

since $\|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|_{2} = \gamma$ as $\boldsymbol{\theta} \in \partial \mathcal{B}_{2}(\boldsymbol{\theta}^{\star}, \gamma)$ and because the Rayleigh quotient of a matrix is bounded below by the smallest eigenvalue of that matrix so that

$$\frac{(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \boldsymbol{H}(\dot{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})}{(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} (\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})} \ \geq \ \lambda_{\min}(\boldsymbol{H}(\dot{\boldsymbol{\theta}})) \ \geq \ \inf_{\boldsymbol{\theta} \in \mathfrak{B}_{2}(\boldsymbol{\theta}^{\star}, \epsilon)} \ \lambda_{\min}(\boldsymbol{H}(\boldsymbol{\theta})),$$

where $\lambda_{\min}(\boldsymbol{H}(\dot{\boldsymbol{\theta}}))$ is the smallest eigenvalue of $\boldsymbol{H}(\dot{\boldsymbol{\theta}})$, noting that

$$\|\dot{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2} = \|t\,\boldsymbol{\theta} + (1-t)\,\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}^{\star}\|_{2} = t\,\|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|_{2} \leq \epsilon,$$

since $t \in [0, 1]$. Lemma 1 showed that

$$\lambda_{\min}(\boldsymbol{H}(\boldsymbol{\theta})) = \lambda_{\min}(I(\boldsymbol{\theta})) \mathbb{E} \|\boldsymbol{Y}\|_{1},$$

which in turn implies

$$\inf_{\boldsymbol{\theta}\in\mathcal{B}_2(\boldsymbol{\theta}^{\star},\epsilon)} \lambda_{\min}(\boldsymbol{H}(\boldsymbol{\theta})) = \widetilde{\lambda}_{\min}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_1,$$

where

$$\widetilde{\lambda}_{\min}^{\epsilon} \coloneqq \inf_{\boldsymbol{\theta} \in \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \epsilon)} \lambda_{\min}(I(\boldsymbol{\theta})),$$

with $I(\boldsymbol{\theta})$ defined in Lemma 1. Hence, for $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \gamma) \ (\gamma \in (0, \epsilon))$,

$$(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \mathbb{E} \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) \geq \gamma^{2} \widetilde{\lambda}_{\min}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_{1} \geq 0.$$
(17)

We next turn to showing

$$\mathbb{P}\left(\inf_{\boldsymbol{\theta}\in\mathfrak{B}_{2}(\boldsymbol{\theta}^{\star},\gamma)}(\boldsymbol{\theta}-\boldsymbol{\theta}^{\star})^{\top}\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y})\geq 0\right) \geq 1-\exp\left(-2p\right)-(\mathbb{E}\|\boldsymbol{Y}\|_{1})^{-1},$$

by showing that the event

$$\sup_{\boldsymbol{\theta}\in \mathfrak{B}_{2}(\boldsymbol{\theta}^{\star},\gamma)} \left| (\boldsymbol{\theta}-\boldsymbol{\theta}^{\star})^{\top} \left(\mathbb{E} \, \nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) - \nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) \right) \right| \ < \ \gamma^{2} \ \widetilde{\lambda}_{\min}^{\epsilon} \ \mathbb{E} \, \|\boldsymbol{Y}\|_{1}$$

occurs with probability at least $1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$. This will in turn imply that the event that $\| \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star} \|_2 \leq \gamma$ will happen with probability at least $1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$. Applying the Cauchy-Schwarz inequality and utilizing standard vector norm inequalities, for $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \gamma)$, we have

$$\begin{split} &|(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})^{\top} \left(\mathbb{E} \, \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) - \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) \right)| \\ &\leq \|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|_{2} \, \|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) - \mathbb{E} \, \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y})\|_{2} \\ &= \gamma \, \|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y}) - \mathbb{E} \, \nabla_{\boldsymbol{\theta}}(\boldsymbol{X}, \boldsymbol{Y})\|_{2}. \end{split}$$

Therefore, it suffices to demonstrate, for all $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \gamma)$,

$$\mathbb{P}\left(\|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) - \mathbb{E}\,\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y})\|_{2} < \gamma\,\widetilde{\lambda}_{\min}^{\epsilon}\,\mathbb{E}\,\|\boldsymbol{Y}\|_{1}\right)$$

is bounded below by

$$1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}.$$

For ease of presentation, we define $\mathcal{D}_{N,\gamma,p}$ to be the event

$$\|
abla_{oldsymbol{ heta}}(oldsymbol{X},oldsymbol{Y}) - \mathbb{E} \,
abla_{oldsymbol{ heta}}(oldsymbol{X},oldsymbol{Y})\|_2 \geq \gamma \, \widetilde{\lambda}_{\min}^{\epsilon} \, \mathbb{E} \, \|oldsymbol{Y}\|_1.$$

Applying Lemma 2, the probability $\mathbb{P}(\mathcal{D}_{N,\gamma,p})$ is bounded above by

$$\exp\left(-\frac{(\gamma \,\widetilde{\lambda}_{\min}^{\epsilon} \,\mathbb{E}\,\|\boldsymbol{Y}\|_{1})^{2}}{36 \,\widetilde{\lambda}_{\max}^{\star} \,(\mathbb{E}\,\|\boldsymbol{Y}\|_{1} + [D_{g}]^{+}) + 2 \sqrt{p} \,\gamma \,\widetilde{\lambda}_{\min}^{\epsilon} \,\mathbb{E}\,\|\boldsymbol{Y}\|_{1}} + p \,\log 5\right) + \frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_{1}},\quad(18)$$

recalling $\widetilde{\lambda}_{\max}^{\star} = \lambda_{\max}(I(\boldsymbol{\theta}^{\star})), \ [D_g]^+ \coloneqq \max\{0, D_g\}, \text{ and }$

$$D_g \coloneqq \sum_{\{i,j\} \prec \{v,w\} \subset \mathbb{N}} \mathbb{C}(Y_{i,j}, Y_{v,w}),$$

where $\{i,j\} \prec \{v,w\}$ implies the sum is taken with respect to the lexicographical ordering of pairs of nodes. Choose

$$\gamma \ = \ \beta \sqrt{\frac{p\,\widetilde{\lambda}_{\max}^{\star}}{\mathbb{E}\|\boldsymbol{Y}\|_1}} \; \frac{1}{\widetilde{\lambda}_{\min}^{\epsilon}},$$

where $\beta>0$ is a positive constant independent of N and p whose value will be determined later. If

$$\lim_{N \to \infty} \beta \sqrt{\frac{p \widetilde{\lambda}_{\max}^{\star}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}}} \frac{1}{\widetilde{\lambda}_{\min}^{\epsilon}} = 0,$$

then for N sufficiently large, we will have $\gamma < \epsilon$, which ensures ϵ may be chosen independent of N and p. While ϵ can be chosen independent of N and p, note that p is expected to be a function of N and thus $\tilde{\lambda}_{\min}^{\epsilon}$ will not (in general) be independent of N, possibly holding implications for how fast p may grow with N for certain θ^{\star} and ϵ . This choice of γ in turn implies that the first term of the exponent in (18) becomes

$$\exp\left(-\frac{\beta^2 \,\widetilde{\lambda}_{\max}^{\star} \mathbb{E} \,\|\boldsymbol{Y}\|_1 \,p}{36 \,\widetilde{\lambda}_{\max}^{\star} \left(\mathbb{E} \,\|\boldsymbol{Y}\|_1 + [D_g]^+\right) + 2 \,\beta \, p \,\sqrt{\mathbb{E} \,\|\boldsymbol{Y}\|_1 \,\widetilde{\lambda}_{\max}^{\star}}}\right).$$
(19)

Canceling $\mathbb{E} \| \boldsymbol{Y} \|_1$ in (19) gives

$$\exp\left(-\frac{\beta^2 \,\widetilde{\lambda}^{\star}_{\max} \, p}{36 \,\widetilde{\lambda}^{\star}_{\max} \, \left(1 + \frac{[D_g]^+}{\mathbb{E} \, \|\boldsymbol{Y}\|_1}\right) \, + \, 2 \,\beta \, p \, \sqrt{\frac{\widetilde{\lambda}^{\star}_{\max}}{\mathbb{E} \, \|\boldsymbol{Y}\|_1}}\right)}.$$

By Assumption 2,

$$p \leq \sqrt{\widetilde{\lambda}_{\max}^{\star} \mathbb{E} \| \boldsymbol{Y} \|_{1}} \left(1 + \frac{[D_{g}]^{+}}{\mathbb{E} \| \boldsymbol{Y} \|_{1}} \right),$$

and by Assumption 1,

$$\frac{[D_g]^+}{\mathbb{E} \|\boldsymbol{Y}\|_1} \leq C_0$$

where $C_0 > 0$ is a positive constant independent of N and p, we have

$$p\sqrt{\frac{\widetilde{\lambda}_{\max}^{\star}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}}} \leq \widetilde{\lambda}_{\max}^{\star} \left(1 + \frac{[D_{g}]^{+}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}}\right) \leq \widetilde{\lambda}_{\max}^{\star} \left(1 + C_{0}\right).$$

As a result, the upper bound in (18) reduces to

$$\mathbb{P}\left(\mathcal{D}_{N,\gamma,p}\right) \leq \exp\left(\left(\frac{-\beta^2}{36(1+C_0)+2(1+C_0)\beta}+\log 5\right)p\right) + \frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_1}.$$
 (20)

To obtain the desired convergence rate, require

$$\frac{-\beta^2}{18\,C+C\,\beta} + \log 5 = -2,\tag{21}$$

where $C = 2(1 + C_0)$. The constant β can be solved by using the quadratic formula and the positive root is given by

$$\beta = \frac{C \log 5 - 2C + \sqrt{(C \log 5 - 2C)^2 + 72(C \log 5 - 2C)}}{2}$$

which ensures $\mathbb{P}(\mathcal{D}_{N,\gamma,p}) \leq \exp(-2p) + (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$. We have thus shown, for all $\boldsymbol{\theta} \in \partial \mathcal{B}_2(\boldsymbol{\theta}^*, \gamma)$, that

$$\mathbb{P}\left(\|\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y}) - \mathbb{E}\nabla_{\boldsymbol{\theta}}(\boldsymbol{X},\boldsymbol{Y})\|_{2} \leq \gamma \,\widetilde{\lambda}_{\min}^{\epsilon} \,\mathbb{E}\,\|\boldsymbol{Y}\|_{1}\right)$$

is bounded below by

$$1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1},$$

under the above conditions. As a result, there exists $N_0 \geq 3$ such that, for all $N \geq N_0$ and with probability at least $1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$, the set $\widehat{\boldsymbol{\Theta}}$ is non-empty and the unique element of the set $\widehat{\boldsymbol{\theta}} \in \widehat{\boldsymbol{\Theta}}$ satisfies (uniqueness following from minimality, as discussed in Section 3)

$$\|\widehat{oldsymbol{ heta}} - oldsymbol{ heta}^{\star}\|_2 \le C rac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{rac{p}{\mathbb{E}\|oldsymbol{Y}\|_1}}.$$

D.1. Proof of Corollary 1

We prove Corollary 1 from Section 3. Under the same assumptions as Theorem 1 and in the case that the parameter dimension p is fixed, the proof of Corollary 1 remains unchanged from that of Theorem 1 except that the exponent in equation (20) scales with N as opposed to p in Theorem 1. Following the same notations in the proof of Theorem 1, we rewrite equation (21) as

$$\frac{-\beta^2}{18\,C+C\,\beta} + \log 5 = -\eta\,(N), \tag{22}$$

where $\eta : \mathbb{N}^+ \to \mathbb{R}^+$ is an increasing function of N. For the ease of notation, we write η_N instead of $\eta(N)$ in the rest of the proof. For any $\alpha_N \in (2 (\mathbb{E} \| \mathbf{Y} \|_1)^{-1}, 1/2)$, let

$$\exp(-\eta_N p) = \frac{\alpha_N}{2}.$$

Note that as N goes to infinity, α_N is allowed to approach 0 through the increasing of η_N . Then the upper bound of $\mathbb{P}(\mathcal{D}_{N,\gamma,p})$ given in (20) becomes

$$\mathbb{P}\left(\mathcal{D}_{N,\gamma,p}\right) \leq \frac{\alpha_N}{2} + \frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_1} \leq \alpha_N,$$

where the last inequality follows from $\alpha_N \geq 2 (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$. To obtain the desired result, solve the positive root of β in terms of η_N from (22):

$$\beta = \frac{C \log 5 + \eta_N C + \sqrt{(C \log 5 + \eta_N C)^2 + 72 (C \log 5 + \eta_N C)}}{2},$$

and write η_N in terms of α_N , where $\alpha_N \to 0$ and $\eta_N \to \infty$ as $N \to \infty$:

$$\eta_N = -\frac{\log\left(\alpha_N / 2\right)}{p}.$$

Let

$$A_1 = C \log 5, \qquad A_2 = \frac{C}{p}.$$

Then for $\alpha_N \in (2 (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}, 1/2),$

$$\beta = \frac{A_1 - A_2 \log(\alpha_N/2) + \sqrt{(A_1 - A_2 \log(\alpha_N/2))^2 + 72 (A_1 - A_2 \log(\alpha_N/2))}}{2}$$

$$= \frac{\log\left(\frac{\alpha_N}{2}\right) \left(\frac{A_1}{\log(\alpha_N/2)} - A_2 + \sqrt{\left(\frac{A_1}{\log(\alpha_N/2)} - A_2\right)^2 + 72 \left(\frac{A_1}{(\log(\alpha_N/2))^2} - \frac{A_2}{\log(\alpha_N/2)}\right)}\right)}{2}$$

$$\leq \frac{\log\left(\frac{\alpha_N}{2}\right) \left(\frac{A_1}{\log 0.25} - A_2 + \sqrt{\left(\frac{A_1}{\log 0.25} - A_2\right)^2 + 72 \left(\frac{A_1}{(\log(0.25))^2} - \frac{A_2}{\log(0.25)}\right)}\right)}{2}$$

$$= A \left|\log\left(\frac{\alpha_N}{2}\right)\right|,$$

where

$$A = \frac{\left|\frac{A_1}{\log 0.25} - A_2 + \sqrt{\left(\frac{A_1}{\log 0.25} - A_2\right)^2 + 72\left(\frac{A_1}{(\log (0.25))^2} - \frac{A_2}{\log (0.25)}\right)}\right|}{2}.$$

As a result, when p is fixed, we showed that for $\alpha_N \in (2 (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}, 1/2)$, with probability at least $1 - \alpha_N$,

$$\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2} \leq A' |\log(\alpha_N/2)| \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{1}{\mathbb{E}\|\boldsymbol{Y}\|_{1}}},$$

where $A' = A\sqrt{p}$ is a positive constant independent of N.

Appendix E: Proof of Theorem 2 and Corollary 2

We prove Theorem 2 and Corollary 2 from Section 3 in one chapter. We first use Fano's method outlined in Chapter 15.3 of Wainwright [43] and the Kullback-Leibler divergence to derive the lower bound of the minimax risk for multilayer network models specified in (1). Let $\epsilon > 0$ be fixed and consider $\gamma \in (0, \epsilon)$. For $M \ge 2$ and some $\delta > 0$, let $\{\theta_1, \ldots, \theta_M\} \subset \mathcal{B}_2(\theta^*, \gamma)$ be a 2δ -separated set. We then have $\|\theta_i - \theta_j\|_2 \ge 2\delta$ for any pair $\{i, j\} \subseteq \{1, \ldots, M\}$. Define the Kullback-Leibler divergence of θ_i and θ_j by

$$\mathrm{KL}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) \coloneqq \sum_{\boldsymbol{x} \in \mathbb{X}} \varphi_{\boldsymbol{\theta}_i}(\boldsymbol{x}) \log \frac{\varphi_{\boldsymbol{\theta}_i}(\boldsymbol{x})}{\varphi_{\boldsymbol{\theta}_j}(\boldsymbol{x})}, \qquad \{i, j\} \subseteq \{1, \dots, M\},$$

where $\varphi_{\theta}(x)$ belongs to a minimal exponential family defined in Proposition 1:

$$\varphi_{\boldsymbol{\theta}}(\boldsymbol{x}) \hspace{2mm} \coloneqq \hspace{2mm} \mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} \mid \boldsymbol{Y} = \boldsymbol{y}) \hspace{2mm} = \hspace{2mm} \exp{(\log{f(\boldsymbol{x}, \boldsymbol{\theta})} + \log{\psi(\boldsymbol{\theta}, \boldsymbol{y})})},$$

recalling $f(\boldsymbol{x}, \boldsymbol{\theta})$ and $\psi(\boldsymbol{\theta}, \boldsymbol{y})$ follow the same form of (1). For $\boldsymbol{\theta} \in \mathbb{R}^p$, denote by $\boldsymbol{s}(\boldsymbol{X}) \in \mathbb{R}^p$ the sufficient statistic vector of the exponential family $\varphi_{\boldsymbol{\theta}}(\boldsymbol{x})$. Then the Kullback-Leibler divergence can be written as

$$\begin{aligned} \mathrm{KL}(\boldsymbol{\theta}_{i},\boldsymbol{\theta}_{j}) &= \sum_{\boldsymbol{x}\in\mathbb{X}} \varphi_{\boldsymbol{\theta}_{i}}(\boldsymbol{x}) \left[\langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, \, \boldsymbol{s}(\boldsymbol{x}) \rangle + \log \, \psi(\boldsymbol{\theta}_{i},\boldsymbol{y}) - \log \, \psi(\boldsymbol{\theta}_{j},\boldsymbol{y}) \right] \\ &= \mathbb{E}_{\boldsymbol{\theta}_{i}} \left\langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, \, \boldsymbol{s}(\boldsymbol{X}) \right\rangle + \log \, \psi(\boldsymbol{\theta}_{i},\boldsymbol{y}) - \log \, \psi(\boldsymbol{\theta}_{j},\boldsymbol{y}) \\ &= \left\langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, \, \boldsymbol{\mu}(\boldsymbol{\theta}_{i}) \right\rangle + \log \, \psi(\boldsymbol{\theta}_{i},\boldsymbol{y}) - \log \, \psi(\boldsymbol{\theta}_{j},\boldsymbol{y}), \end{aligned}$$

$$\end{aligned}$$

where $\mu(\theta) \coloneqq \mathbb{E}_{\theta} s(X)$ is the mean-value parameter map of the exponential family. By Corollary 2.3 of Brown [7],

$$\log \psi(\boldsymbol{\theta}_{j}) = \log \psi(\boldsymbol{\theta}_{i}) + \langle \boldsymbol{\theta}_{j} - \boldsymbol{\theta}_{i}, -\boldsymbol{\mu}(\boldsymbol{\theta}_{i}) \rangle - \frac{1}{2} \langle \boldsymbol{\theta}_{j} - \boldsymbol{\theta}_{i}, I_{\boldsymbol{X}}(\dot{\boldsymbol{\theta}}) (\boldsymbol{\theta}_{j} - \boldsymbol{\theta}_{i}) \rangle$$

$$= \log \psi(\boldsymbol{\theta}_{i}) + \langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, \boldsymbol{\mu}(\boldsymbol{\theta}_{i}) \rangle - \frac{1}{2} \langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, I_{\boldsymbol{X}}(\dot{\boldsymbol{\theta}}) (\boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}) \rangle, \qquad (24)$$

where $\dot{\boldsymbol{\theta}} = t\boldsymbol{\theta}_i + (1-t)\boldsymbol{\theta}_j$ for some $t \in (0,1)$, and $I_{\boldsymbol{X}}(\dot{\boldsymbol{\theta}})$ is the Fisher information matrix at $\dot{\boldsymbol{\theta}}$ for $\boldsymbol{X} \in \mathbb{X}$. For a fixed $\epsilon > 0$ such that $\gamma \in (0,\epsilon)$ and $\{\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_M\} \subset \mathcal{B}_2(\boldsymbol{\theta}^{\star}, \gamma)$, define

$$\widetilde{\lambda}_{\max}^{\epsilon} \quad \coloneqq \quad \sup_{\boldsymbol{\theta} \in \mathcal{B}_{2}(\boldsymbol{\theta}^{\star}, \epsilon)} \ \frac{\lambda_{\max}(I_{\boldsymbol{X}}(\boldsymbol{\theta}))}{\mathbb{E} \, \|\boldsymbol{Y}\|_{1}} \quad = \quad \sup_{\boldsymbol{\theta} \in \mathcal{B}_{2}(\boldsymbol{\theta}^{\star}, \epsilon)} \lambda_{\max} \, (I(\boldsymbol{\theta})),$$

where $\lambda_{\max}(\mathbf{A})$ is the maximum eigenvalue of matrix \mathbf{A} and $I(\boldsymbol{\theta})$ is the Fisher information matrix for an activated dyad defined in Lemma 1. Combining (23) and (24) and using the standard matrix norm inequality and the triangle inequality, we have

$$\begin{split} \operatorname{KL}(\boldsymbol{\theta}_{i},\boldsymbol{\theta}_{j}) &= \frac{1}{2} \left\langle \boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}, \ I_{\boldsymbol{X}}(\dot{\boldsymbol{\theta}}) \left(\boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}\right) \right\rangle \\ &\leq \frac{1}{2} \operatorname{\mathbb{E}} \|\boldsymbol{Y}\|_{1} \widetilde{\lambda}_{\max}^{\epsilon} \|\boldsymbol{\theta}_{i} - \boldsymbol{\theta}_{j}\|_{2}^{2} \\ &\leq \frac{1}{2} \operatorname{\mathbb{E}} \|\boldsymbol{Y}\|_{1} \widetilde{\lambda}_{\max}^{\epsilon} \left(\|\boldsymbol{\theta}_{i} - \boldsymbol{\theta}^{\star}\|_{2} + \|\boldsymbol{\theta}_{j} - \boldsymbol{\theta}^{\star}\|_{2}\right)^{2} \\ &\leq 2 \epsilon^{2} \operatorname{\mathbb{E}} \|\boldsymbol{Y}\|_{1} \widetilde{\lambda}_{\max}^{\epsilon}. \end{split}$$

Note that the size M of the largest possible 2δ -separated set $\{\theta_1, \ldots, \theta_M\} \subset \mathcal{B}_2(\theta^*, \gamma) \subset \mathbb{R}^p$ is the packing number of $\mathcal{B}_2(\theta^*, \gamma)$. By Lemma 4.2.8 and Corollary 4.2.13 of Vershynin [42], we have

$$M \geq \left(\frac{\gamma}{2\,\delta}\right)^p,$$

and

$$\log M \geq p \log \left(\frac{\gamma}{2\delta}\right).$$

By Proposition 15.12 of Wainwright [43], the minimax risk \mathcal{R}_N has the lower bound

$$\mathcal{R}_N \geq \delta \left[1 - \frac{\mathcal{F} + \log 2}{\log M} \right],$$

where

$$\mathcal{F} := \max_{\{i,j\} \subseteq \{1,\dots,M\}} \operatorname{KL}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j).$$

Since $\gamma \in (0, \epsilon)$, the lower bound for \mathcal{R}_N can be written as

$$\mathcal{R}_N \geq \delta \left[1 - \frac{2 \gamma^2 \mathbb{E} \| \boldsymbol{Y} \|_1 \widetilde{\lambda}_{\max}^{\epsilon} + \log 2}{p \log (\gamma/2 \delta)} \right].$$

To obtain the desired lower bound $\mathcal{R}_N \geq \delta/2$, we need

$$\frac{2\gamma^2 \mathbb{E} \, \|\boldsymbol{Y}\|_1 \, \widetilde{\lambda}_{\max}^\epsilon + \log 2}{p \, \log \, (\gamma/2 \, \delta)} \quad \leq \quad \frac{1}{2},$$

which implies

$$\frac{4\gamma^2 \mathbb{E} \|\boldsymbol{Y}\|_1 \,\widetilde{\lambda}_{\max}^{\epsilon}}{p} + \frac{2\log\left(2\right)}{p} \leq \log(\gamma/2) - \log\left(\delta\right).$$

Exponentiating both sides we have

$$\exp\left(\frac{4\gamma^2 \mathbb{E} \|\boldsymbol{Y}\|_1 \widetilde{\lambda}_{\max}^{\epsilon}}{p} + \frac{2\log\left(2\right)}{p}\right) \leq \frac{\gamma/2}{\delta}.$$

This leads us to the following inequality

$$\delta \leq \frac{\gamma}{2} \exp\left(-\frac{4\gamma^2 \mathbb{E} \|\boldsymbol{Y}\|_1 \,\widetilde{\lambda}_{\max}^{\epsilon}}{p} - \frac{2\log\left(2\right)}{p}\right).$$

Choosing

$$\gamma = 2 C \sqrt{\frac{p}{\widetilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_{1}}}$$

for some C > 0, we obtain the bound

$$\delta \leq C \exp\left(-16 C^2 - \frac{2 \log(2)}{p}\right) \sqrt{\frac{p}{\widetilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_1}}.$$
 (25)

As long as $p = O(\widetilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \| \boldsymbol{Y} \|_1)$, we can choose C to ensure $\gamma \in (0, \epsilon)$. Finally, for all $\delta > 0$ satisfying (25), we have \mathcal{R}_N lower bounded by

$$\mathcal{R}_N \geq \frac{\delta}{2}.$$

Note that as $p \ge 1$,

$$\exp\left(-16 C^2 - \frac{2 \log(2)}{p}\right) \ge \exp\left(-16 C^2 - 2 \log(2)\right),$$

we may choose

$$\delta = C \exp\left(-16 C^2 - 2 \log\left(2\right)\right) \sqrt{\frac{p}{\tilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_1}}$$
$$= A' \sqrt{\frac{p}{\tilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_1}},$$

where $A' = C \exp(-16 C^2 - 2 \log(2))$. Then we obtain the desired lower bound for the minimax risk

$$\mathcal{R}_N \geq \frac{\delta}{2} = \frac{A'}{2} \sqrt{\frac{p}{\widetilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\boldsymbol{Y}\|_1}}.$$
 (26)

Next, we show the lower bound in (26) matches with the upper bound of the ℓ_2 -error of the maximum likelihood estimator $\hat{\theta}$ provided in Theorem 1. Let A = A'/2 be an unknown constant independent of N and p. We have

$$\begin{aligned} \mathcal{R}_{N} &\geq A \sqrt{\frac{p}{\widetilde{\lambda}_{\max}^{\epsilon} \mathbb{E} \|\mathbf{Y}\|_{1}}} \\ &= A \sqrt{\frac{\widetilde{\lambda}_{\max}^{\star}}{\widetilde{\lambda}_{\max}^{\star}}} \left(\frac{\widetilde{\lambda}_{\min}^{\epsilon}}{\widetilde{\lambda}_{\min}^{\epsilon}}\right) \frac{1}{\sqrt{\widetilde{\lambda}_{\max}^{\epsilon}}} \sqrt{\frac{p}{\mathbb{E} \|\mathbf{Y}\|_{1}}} \\ &= A \frac{1}{\sqrt{\widetilde{\lambda}_{\max}^{\epsilon}}} \frac{\widetilde{\lambda}_{\min}^{\epsilon}}{\sqrt{\widetilde{\lambda}_{\max}^{\star}}} \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{p}{\mathbb{E} \|\mathbf{Y}\|_{1}}} \\ &\geq A \left(\frac{\widetilde{\lambda}_{\min}^{\epsilon}}{\widetilde{\lambda}_{\max}^{\epsilon}}\right) \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{p}{\mathbb{E} \|\mathbf{Y}\|_{1}}}, \end{aligned}$$

where the last inequality holds because $\tilde{\lambda}_{\max}^{\epsilon} \geq \tilde{\lambda}_{\max}^{\star}$. Under the assumption that

$$\widetilde{\lambda}_{\max}^{\epsilon} = O\left(\widetilde{\lambda}_{\min}^{\epsilon}\right),\,$$

we showed that the lower bound of the minimax risk \mathcal{R}_N and the upper bound of the ℓ_2 -error of the maximum likelihood estimator presented in Theorem 1 match up to an unknown constant independent of N and p.

Appendix F: Proposition 2 and proof

In order to establish a bound on the error of the multivariate normal approximation for estimators of data-generating parameters, we first establish an error bound on the multivariate normal approximation of a standardization of the sufficient statistic vector s(X) of the exponential family distribution of X given Y, derived in Lemma 3, in Proposition 2 using a Lyapunov type bound presented in Raič [33]. Proposition 2 provides the basis for our normality proof for estimators which we present in Theorem 3.

Proposition 2 Consider a separable multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ nodes and $K \ge 1$ layers. Denote by $\mathbf{s}(\mathbf{X}) \in \mathbb{R}^p$ the sufficient statistic vector of the exponential family $\mathbb{P}(\mathbf{X} = \mathbf{x} | \mathbf{Y} = \mathbf{y})$ as defined in Lemma 3. Let $\mathbb{E}^{\mathbf{Y}}$ be the random conditional expectation operator for the distribution of \mathbf{X} conditional on \mathbf{Y} , and define

$$egin{aligned} oldsymbol{S}_{\mathcal{N}} &\coloneqq & (I(oldsymbol{ heta}^{\star}) \, \|oldsymbol{Y}\|_1)^{-1/2} \left(oldsymbol{s}(oldsymbol{X}) - \mathbb{E}^{oldsymbol{Y}} oldsymbol{s}(oldsymbol{X}))
ight) \ &= & \sum_{\{i,j\} \subset \mathcal{N}} \left(I(oldsymbol{ heta}^{\star}) \, \|oldsymbol{Y}\|_1
ight)^{-1/2} \left(oldsymbol{s}_{i,j}(oldsymbol{X}) - \mathbb{E}^{oldsymbol{Y}} oldsymbol{s}_{i,j}(oldsymbol{X})). \end{aligned}$$

For any measurable convex set $\mathcal{A} \subset \mathbb{R}^p$,

$$\left| \mathbb{P}(\boldsymbol{S}_{\mathcal{N}} \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A}) \right| \leq \frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}} \sqrt{\frac{p^{7/2}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{4}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{8 [D_{g}]^{+}}{(\mathbb{E} \|\boldsymbol{Y}\|_{1})^{2}}},$$

where Φ is the standard multivariate normal measure and $\mathbf{Z} \sim MvtNorm(\mathbf{0}_p, \mathbf{I}_p)$, where $\mathbf{0}_p$ is the p-dimensional vector of zeros and \mathbf{I}_p is the $p \times p$ identity matrix.

Before we prove Proposition 2, we introduce a Lyapunov type bound in Lemma 4 provided by Theorem 1 of Raic [33].

Lemma 4. Consider a sequence of $n \ge 1$ independent random vectors $\mathbf{W}_i \in \mathbb{R}^p$. Assume that $\mathbb{E} \mathbf{W}_i = \mathbf{0}_p$ and $\sum_{i=1}^n \mathbb{V} \mathbf{W}_i = \mathbf{I}_p$ where $\mathbf{0}_p$ is the p-dimensional vector of zeros and \mathbf{I}_p is the $p \times p$ identity matrix. Define

$$oldsymbol{S}_n \hspace{.1in} = \hspace{.1in} \sum_{i=1}^n oldsymbol{W}_i$$

and let \mathbf{Z} be the standard multivariate normal random variable, i.e., $\mathbf{Z} \sim MvtNorm(\mathbf{0}_p, \mathbf{I}_p)$. Then, for all measurable convex sets $\mathcal{A} \subset \mathbb{R}^p$,

$$|\mathbb{P}(\boldsymbol{S}_n \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \leq (42 p^{1/4} + 16) \sum_{i=1}^n \mathbb{E} \|\boldsymbol{W}_i\|_2^3,$$

where Φ is the standard multivariate normal measure.

We now turn to proving Proposition 2.

PROOF OF PROPOSITION 2. By Proposition 1 and Lemma 3, the conditional distribution of the multilayer network X given Y follows an exponential family with sufficient statistic vector that can be decomposed into the sum of conditionally independent dyad-based statistics:

$$oldsymbol{s}(oldsymbol{X}) \hspace{.1in} = \hspace{.1in} \sum_{\{i,j\} \subset \mathcal{N}} oldsymbol{s}_{i,j}(oldsymbol{X}),$$

with the precise formula for $s_{i,j}(X)$ given in Lemma 3. Define

$$\begin{split} \boldsymbol{S}_{\mathcal{N}} &\coloneqq (I(\boldsymbol{\theta}^{\star}) \, \|\boldsymbol{Y}\|_{1})^{-1/2} \left(\boldsymbol{s}(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} \boldsymbol{s}(\boldsymbol{X}) \right) \\ &= \sum_{\{i,j\} \subset \mathcal{N}} (I(\boldsymbol{\theta}^{\star}) \, \|\boldsymbol{Y}\|_{1})^{-1/2} \left(\boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} \boldsymbol{s}_{i,j}(\boldsymbol{X}) \right), \end{split}$$

where $I(\boldsymbol{\theta}^{\star})$ is the Fisher information matrix of an activated dyad $X_{i,j}$ for $\{i, j\} \subset \mathbb{N}$ satisfying $Y_{i,j} = 1$ evaluated at $\boldsymbol{\theta}^{\star}$ per Lemma 1 and where $\mathbb{E}^{\mathbf{Y}}$ is the random conditional expectation operator with respect to the distribution of \boldsymbol{X} conditional on \boldsymbol{Y} . For $\gamma > 0$ satisfying $\gamma < \mathbb{E} \|\boldsymbol{Y}\|_1$, define the event $\mathcal{E}(\gamma)$ by

$$\mathcal{E}(\gamma) \cong \{ \boldsymbol{y} \in \mathbb{Y} : \| \boldsymbol{y} \|_1 \ge \mathbb{E} \| \boldsymbol{Y} \|_1 - \gamma \}.$$

In words, $\mathcal{E}(\gamma)$ is the subset of configurations of the single-layer network \boldsymbol{Y} which have the number of edges equal to at least the expected number of activated dyads $\mathbb{E} \|\boldsymbol{Y}\|_1$ minus $\gamma > 0$. The restrictions placed on γ ensure that $\mathbb{E} \|\boldsymbol{Y}\|_1 - \gamma >$ 0 which implies that $\mathcal{E}(\gamma)$ will not contain the empty graph which has no edges and that $\mathcal{E}(\gamma)$ will contain the complete graph with $\binom{N}{2}$ edges as $\mathbb{E} \|\boldsymbol{Y}\|_1 < \binom{N}{2}$ (strict inequality following from the fact that $g(\boldsymbol{y})$, the marginal probability mass function of \boldsymbol{Y} , is assumed to be strictly positive on \mathbb{Y}). Hence, $\mathbb{P}(\mathcal{E}(\gamma)) > 0$ and $\mathbb{P}(\mathcal{E}(\gamma)^c) > 0$. Let $\mathcal{A} \subset \mathbb{R}^p$ be a measurable convex set. By the law of total probability and the triangle inequality, we have

$$|\mathbb{P}(\boldsymbol{S}_{\mathcal{N}} \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \leq |\mathbb{P}(\boldsymbol{S}_{n} \in \mathcal{A} | \mathcal{E}(\gamma)) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \mathbb{P}(\mathcal{E}(\gamma)) + |\mathbb{P}(\boldsymbol{S}_{n} \in \mathcal{A} | \mathcal{E}^{c}(\gamma)) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \mathbb{P}(\mathcal{E}^{c}(\gamma)) \leq \sup_{\boldsymbol{y} \in \mathcal{E}(\gamma)} |\mathbb{P}(\boldsymbol{S}_{\mathcal{N}} \in \mathcal{A} | \boldsymbol{Y} = \boldsymbol{y}) - \Phi(\boldsymbol{Z} \in \mathcal{A})| + \mathbb{P}(\mathcal{E}^{c}(\gamma)),$$

$$(27)$$

noting $|\mathbb{P}(\boldsymbol{S}_n \in \mathcal{A} | \mathcal{E}^c(\gamma)) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \leq 1$ and $\mathbb{P}(\mathcal{E}(\gamma)) \leq 1$. Taking

$$\boldsymbol{W}_{i,j} = (I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_1)^{-1/2} \left(\boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} \, \boldsymbol{s}_{i,j}(\boldsymbol{X}) \right),$$

we have

$$\mathbb{E}\left[\boldsymbol{W}_{i,j} \,|\, \boldsymbol{Y} = \boldsymbol{y}\right] = 0$$

a result of the tower property of conditional expectation, and

$$\mathbb{V}\left[\sum_{\{i,j\}\subset\mathcal{N}} \boldsymbol{W}_{i,j} \mid \boldsymbol{Y} = \boldsymbol{y}\right] = \boldsymbol{I}_p,$$

which follows from Lemma 1 which establishes that $\mathbb{V}[s_{i,j}(\mathbf{X}) | \mathbf{Y} = \mathbf{y}] = I(\boldsymbol{\theta}^*)$ when $Y_{i,j} = 1$, recalling the form of the Fisher information matrix of exponential families to be the covariance matrix of the sufficient statistic vector [e.g., Proposition 3.10, pp. 32, 41], and due to the fact that $\mathbb{V}[s_{i,j}(\mathbf{X}) | \mathbf{Y} = \mathbf{y}] = \mathbf{0}_{p,p}$ when $Y_{i,j} = 0$. Applying Lemma 4 to the first term of the summation of (27), for any measurable convex set $\mathcal{A} \subset \mathbb{R}^p$,

$$|\mathbb{P}(\boldsymbol{S}_{\mathcal{N}} \in \mathcal{A} \mid \boldsymbol{Y} = \boldsymbol{y}) - \Phi(\boldsymbol{Z} \in \mathcal{A})|$$

is bounded above by

$$(42 p^{1/4} + 16) \sum_{\{i,j\} \subset \mathcal{N}} \mathbb{E} \left[\| \boldsymbol{W}_{i,j} \|_2^3 \, | \, \boldsymbol{Y} = \boldsymbol{y} \right].$$

Given Y = y, using standard matrix and vector norm inequalities,

$$\begin{split} \| \boldsymbol{W}_{i,j} \|_{2} &= \| (I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{y} \|_{1})^{-1/2} \left(\boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E} \, \boldsymbol{s}_{i,j}(\boldsymbol{X}) \right) \|_{2} \\ &\leq \| \boldsymbol{y} \|_{1}^{-1/2} \| \| I(\boldsymbol{\theta}^{\star})^{-1/2} \|_{2} \| \boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E} \, \boldsymbol{s}_{i,j}(\boldsymbol{X}) \|_{2} \\ &\leq (\| \boldsymbol{y} \|_{1} \, \widetilde{\lambda}_{\min}^{\epsilon})^{-1/2} \, p^{1/2} \, y_{i,j}, \end{split}$$

where $\| \cdot \|_2$ denotes the spectral norm of a $p \times p$ matrix. From the proof of Lemma 2, for all $l \in \{1, \ldots, p\}$, we have

$$0 \leq s_{l,i,j}(\boldsymbol{x}) \leq 1, \quad \{i,j\} \subset \mathcal{N},$$

 $\mathbb P\text{-almost}$ surely. Hence,

$$\mathbb{P}(\|\boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}}\boldsymbol{s}_{i,j}(\boldsymbol{X})\|_{\infty} \le y_{i,j} \mid \boldsymbol{Y} = \boldsymbol{y}) = 1,$$

implying (conditional on Y = y)

$$\|\boldsymbol{s}_{i,j}(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} \, \boldsymbol{s}_{i,j}(\boldsymbol{X})\|_2 \le (p \, y_{i,j})^{1/2} = p^{1/2} \, y_{i,j},$$

 \mathbb{P} -almost surely. As a result,

$$\mathbb{E}\left[\|\boldsymbol{W}_{i,j}\|_2^3 \,|\, \boldsymbol{Y} = \boldsymbol{y}\right] \leq (\|\boldsymbol{y}\|_1 \,\widetilde{\lambda}_{\min}^{\epsilon})^{-3/2} \, p^{3/2} \, y_{i,j},$$

noting that $y_{i,j}^3 = y_{i,j} \in \{0,1\}$. Using the fact that $42 p^{1/4} + 16 \le 58 p^{1/4} (p \ge 1)$, we have

$$(42 p^{1/4} + 16) \sum_{\{i,j\} \subset \mathcal{N}} \mathbb{E} \left[\| \boldsymbol{W}_{i,j} \|_{2}^{3} | \boldsymbol{Y} = \boldsymbol{y} \right]$$

$$\leq 58 p^{7/4} \sum_{\{i,j\} \subset \mathcal{N}} y_{i,j} \left(\| \boldsymbol{y} \|_{1} \widetilde{\lambda}_{\min}^{\epsilon} \right)^{-3/2}$$

$$= 58 p^{7/4} \| \boldsymbol{y} \|_{1}^{-1/2} (\widetilde{\lambda}_{\min}^{\epsilon})^{-3/2}$$

$$\leq 58 p^{7/4} \left(\mathbb{E} \| \boldsymbol{Y} \|_{1} - \gamma \right)^{-1/2} (\widetilde{\lambda}_{\min}^{\epsilon})^{-3/2},$$

as the conditioning event $\mathcal{E}(\gamma)$ and choice of γ ensure that $\|\boldsymbol{y}\|_1 \geq \mathbb{E}\|\boldsymbol{Y}\|_1 - \gamma > 0$. We bound the second term in (27) by Chebyshev's inequality using equation (14) in the proof of Lemma 2:

$$\mathbb{P}(\mathcal{E}^{c}(\gamma)) \leq \frac{\mathbb{E} \|\boldsymbol{Y}\|_{1} + 2 [D_{g}]^{+}}{\gamma^{2}}.$$

Taking $\gamma = 2^{-1} \mathbb{E} \| \boldsymbol{Y} \|_1 > 0$, we have

$$\mathbb{P}(\mathcal{E}^{c}(\gamma)) \leq \frac{4}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{8 \left[D_{g}\right]^{+}}{\left(\mathbb{E} \|\boldsymbol{Y}\|_{1}\right)^{2}}.$$

Combining terms, we obtain the bound

$$|\mathbb{P}(\boldsymbol{S}_{\mathcal{N}} \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})| \leq \frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}} \sqrt{\frac{p^{7/2}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{4}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{8 [D_{g}]^{+}}{(\mathbb{E} \|\boldsymbol{Y}\|_{1})^{2}}}.$$

Note that the asymptotic multivariate normality can be established provided

$$\lim_{N \to \infty} \left[\frac{83}{(\tilde{\lambda}_{\min}^{\epsilon})^{3/2}} \sqrt{\frac{p^{7/2}}{\mathbb{E} \|\boldsymbol{Y}\|_{1}}} + \frac{4}{\mathbb{E} \|\boldsymbol{Y}\|_{1}} + \frac{8 [D_{g}]^{+}}{(\mathbb{E} \|\boldsymbol{Y}\|_{1})^{2}} \right] = 0,$$

resulting in the following asymptotic convergence in distribution:

$$\boldsymbol{S}_{\mathcal{N}} \stackrel{D}{\longrightarrow} \boldsymbol{Z} \sim \operatorname{MvtNorm}\left(\boldsymbol{0}, \boldsymbol{I}_{p}\right).$$

Appendix G: Proof of Theorem 3

In order to demonstrate the feasibility of the normal approximation for maximum likelihood estimators $\hat{\theta}$ of θ^{\star} , we first start with a standard Taylor expansion of the negative score equation:

$$-\nabla_{\boldsymbol{\theta}} \,\ell(\widehat{\boldsymbol{\theta}}; \boldsymbol{x}, \boldsymbol{y}) = -\nabla_{\boldsymbol{\theta}} \,\ell(\boldsymbol{\theta}^{\star}; \boldsymbol{x}, \boldsymbol{y}) - \nabla_{\boldsymbol{\theta}}^{2} \,\ell(\boldsymbol{\theta}^{\star}; \boldsymbol{x}, \boldsymbol{y}) \,(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \boldsymbol{R}, \qquad (28)$$

where $\mathbf{R} \in \mathbb{R}^p$ is the vector of remainders given in the Lagrange form. Denoting by R_i , $(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star})_i$, and $(\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}))_i$ the *i*th component of \mathbf{R} , $(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star})$, and the score function $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$, respectively. The remainder term R_i $(i = 1, \ldots, p)$ is given by

$$R_{i} = \sum_{j=1}^{p} \frac{1}{2} \frac{\partial^{2} (\nabla_{\boldsymbol{\theta}} \, \ell(\dot{\boldsymbol{\theta}}_{i}; \boldsymbol{x}, \boldsymbol{y}))_{i}}{\partial \, \theta_{j}^{2}} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star})_{j}^{2} \\ + \sum_{1 \leq j < k \leq p} \frac{\partial^{2} (\nabla_{\boldsymbol{\theta}} \ell(\dot{\boldsymbol{\theta}}_{i}; \boldsymbol{x}, \boldsymbol{y}))_{i}}{\partial \, \theta_{j} \, \partial \, \theta_{k}} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star})_{j} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star})_{k},$$

where $\dot{\boldsymbol{\theta}}_i = t_i \, \widehat{\boldsymbol{\theta}} + (1 - t_i) \, \boldsymbol{\theta}^{\star}$ (for some $t_i \in [0, 1]$). By Proposition 1,

$$\ell(\boldsymbol{ heta}; \boldsymbol{x}, \boldsymbol{y}) = \log \mathbb{P}_{\boldsymbol{ heta}}(\boldsymbol{X} = \boldsymbol{x} \,|\, \boldsymbol{Y} = \boldsymbol{y}) + \log g(\boldsymbol{y})$$

By Lemma 3, the probability mass function $\mathbb{P}_{\theta}(X = x | Y = y)$ belongs to a minimal exponential family with the sufficient statistic vector s(x) given by equation (15) in Lemma 3. We then have,

$$\begin{aligned} -\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) &= -(\boldsymbol{s}(\boldsymbol{x}) - \mathbb{E}_{\boldsymbol{\theta}}^{\boldsymbol{y}} \, \boldsymbol{s}(\boldsymbol{X})) \\ -\nabla_{\boldsymbol{\theta}}^{2} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) &= \mathbb{V}_{\boldsymbol{\theta}}^{\boldsymbol{y}} \, \boldsymbol{s}(\boldsymbol{X}) &= I(\boldsymbol{\theta}^{\star}) \, \|\boldsymbol{y}\|_{1} \end{aligned}$$

where $\mathbb{E}_{\theta}^{\boldsymbol{y}}$ and $\mathbb{V}_{\theta}^{\boldsymbol{y}}$ are the conditional expectation and variance operators, respectively, of the conditional distribution of \boldsymbol{X} given $\boldsymbol{Y} = \boldsymbol{y}$, and by using standard formulas for exponential families [e.g., Proposition 3.8, pp. 29, 41] and the results of Lemma 1. Note $\nabla_{\theta} \ell(\hat{\theta}; \boldsymbol{x}, \boldsymbol{y}) = 0$, as the maximum likelihood estimator $\hat{\boldsymbol{\theta}}$ solves the score equation by definition. We re-arrange (28) and multiply both sides by $(I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_1)^{-1/2}$ to obtain

$$(I(\boldsymbol{\theta}^{\star}) \|\boldsymbol{Y}\|_{1})^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - (I(\boldsymbol{\theta}^{\star}) \|\boldsymbol{Y}\|_{1})^{-1/2} \boldsymbol{R}$$

= $(I(\boldsymbol{\theta}^{\star}) \|\boldsymbol{Y}\|_{1})^{-1/2} (s(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} s(\boldsymbol{X})).$ (29)

Define $\Delta := (I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_1)^{-1/2} \boldsymbol{R}$. Let $\mathcal{A} \subset \mathbb{R}^p$ be any measurable convex subset of \mathbb{R}^p and $\boldsymbol{Z} \sim \text{MvtNorm}(\boldsymbol{0}_p, \boldsymbol{I}_p)$. We are interested in bounding the quantity

$$\left|\mathbb{P}((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \Delta \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})\right|.$$

Then from (29),

$$\mathbb{P}\left((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} \left(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\right) - \Delta \in \mathcal{A}\right)$$

= $\mathbb{P}\left((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{-1/2} \left(s(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} s(\boldsymbol{X})\right) \in \mathcal{A}\right).$

Applying Proposition 2, for all measurable convex sets $\mathcal{A} \subseteq \mathbb{R}^p$,

$$\left|\mathbb{P}\left(\left(I(\boldsymbol{\theta}^{\star}) \|\boldsymbol{Y}\|_{1}\right)^{-1/2} \left(s(\boldsymbol{X}) - \mathbb{E}^{\boldsymbol{Y}} s(\boldsymbol{X})\right) \in \mathcal{A}\right) - \Phi(\boldsymbol{Z} \in \mathcal{A})\right|$$

is bounded above by

$$\frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}}\sqrt{\frac{p^{7/2}}{\mathbb{E}\,\|\boldsymbol{Y}\|_1}} + \frac{4}{\mathbb{E}\,\|\boldsymbol{Y}\|_1} + \frac{8\,\left[D_g\right]^+}{\left(\mathbb{E}\,\|\boldsymbol{Y}\|_1\right)^2}.$$

Hence,

$$\left|\mathbb{P}((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \Delta \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})\right|$$

is bounded above by

$$\frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}}\sqrt{\frac{p^{7/2}}{\mathbb{E}\,\|\boldsymbol{Y}\|_1}} + \frac{4}{\mathbb{E}\,\|\boldsymbol{Y}\|_1} + \frac{8\,[D_g]^+}{\left(\mathbb{E}\,\|\boldsymbol{Y}\|_1\right)^2}.$$

24

We lastly handle the term Δ by showing that $\|\Delta\|_2$ is small with high probability. We first use standard vector/matrix norm inequalities to bound

$$\|\Delta\|_{2} = \|(I(\boldsymbol{\theta}^{\star}) \,\|\boldsymbol{Y}\|_{1})^{-1/2} \,\boldsymbol{R}\|_{2} \leq \frac{\|I(\boldsymbol{\theta}^{\star})^{-1/2}\|_{2}}{\sqrt{\|\boldsymbol{Y}\|_{1}}} \,\|\boldsymbol{R}\|_{2} \leq \frac{\|\boldsymbol{R}\|_{2}}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}}\|\boldsymbol{Y}\|_{1}},$$

noting that the spectral norm $|||I(\theta^*)^{-1/2}|||_2$ is equal to the largest eigenvalue of $I(\theta^*)^{-1/2}$ which will be the reciprocal of the smallest eigenvalue of $I(\theta^*)^{1/2}$, which is bounded below by $\sqrt{\tilde{\lambda}_{\min}^{\epsilon}}$. Using a standard result from the Taylor theorem for functions with multiple variables, if for each $i = 1, \ldots, p$, there exists constants $M_i > 0$ such that

$$\sup_{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_1 \le \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star\|_1} \left| \frac{\partial^2 (\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}))_i}{\partial \, \theta_j \, \partial \, \theta_k} \right| \le M_i, \qquad 1 \le j \le k \le p,$$

then the Lagrange remainder is bounded above by

$$R_i \leq \frac{M_i}{2} \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star\|_1^2$$

on the set $\{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_1 \leq \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star\|_2\}$. By Lemma 5, conditional on $\boldsymbol{Y} = \boldsymbol{y}$, we have, for all $i = 1, \ldots, p$, the bound $M_i \leq 2 \|\boldsymbol{y}\|_1$. Hence,

$$\begin{split} \|\Delta\|_{2} &\leq \frac{1}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}} \|\boldsymbol{y}\|_{1}} \sqrt{\sum_{i=1}^{p} R_{i}^{2}} \leq \frac{1}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}} \|\boldsymbol{y}\|_{1}} \sqrt{\sum_{i=1}^{p} \|\boldsymbol{y}\|_{1}^{2} \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{1}^{4}} \\ &\leq \frac{1}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}} \|\boldsymbol{y}\|_{1}} \sqrt{p \|\boldsymbol{y}\|_{1}^{2} \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{1}^{4}} \leq \frac{\sqrt{p} \|\boldsymbol{y}\|_{1} \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{1}^{2}}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}} \|\boldsymbol{y}\|_{1}} \tag{30} \\ &\leq \frac{\sqrt{p} \sqrt{\|\boldsymbol{y}\|_{1}} p \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2}^{2}}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}}} \leq \frac{p^{3/2} \sqrt{\|\boldsymbol{y}\|_{1}} \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2}^{2}}{\sqrt{\widetilde{\lambda}_{\min}^{\epsilon}}}. \end{split}$$

By Chebyshev's inequality—as in the proof of Lemma 2—we can establish that

$$\mathbb{P}\left(\left|\|\boldsymbol{Y}\|_{1} - \mathbb{E}\|\boldsymbol{Y}\|_{1}\right| > \frac{1}{2}\mathbb{E}\|\boldsymbol{Y}\|_{1}\right) \leq \frac{4}{\mathbb{E}\|\boldsymbol{Y}\|_{1}} + \frac{8[D_{g}]^{+}}{(\mathbb{E}\|\boldsymbol{Y}\|_{1})^{2}}.$$
 (31)

Under Assumptions 1, 2 and 3, Theorem 1 established that there exist constants C > 0 and $N_0 \ge 3$ such that, for all $N \ge N_0$, the event

$$\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}\|_{2} \leq C \frac{\sqrt{\widetilde{\lambda}_{\max}^{\star}}}{\widetilde{\lambda}_{\min}^{\epsilon}} \sqrt{\frac{p}{\mathbb{E}\|\boldsymbol{Y}\|_{1}}}$$
(32)

occurs with probability at least $1 - \exp(-2p) - (\mathbb{E} \| \boldsymbol{Y} \|_1)^{-1}$. Define \mathcal{E}_1 to be the event

$$\|\|oldsymbol{Y}\|_1 - \mathbb{E}\,\|oldsymbol{Y}\|_1| \leq rac{1}{2}\,\mathbb{E}\,\|oldsymbol{Y}\|_1$$

and \mathcal{E}_2 to be the event in (32), and define \mathcal{R} to be the corresponding values of Δ in the event $(\mathbf{X}, \mathbf{Y}) \in \mathcal{E}_1 \cap \mathcal{E}_2$, under which we have the bound

$$\begin{split} \|\Delta\|_{2} &\leq \frac{p^{3/2} \sqrt{\|\boldsymbol{y}\|_{1}}}{\sqrt{\tilde{\lambda}_{\min}^{\epsilon}}} C^{2} \frac{\tilde{\lambda}_{\max}^{\star}}{(\tilde{\lambda}_{\min}^{\epsilon})^{2}} \frac{p}{\mathbb{E}\|\boldsymbol{Y}\|_{1}} \\ &\leq \frac{C^{2} p^{5/2} \sqrt{2\mathbb{E}} \|\boldsymbol{Y}\|_{1}}{\mathbb{E}\|\boldsymbol{Y}\|_{1}} \frac{\tilde{\lambda}_{\max}^{\star}}{(\tilde{\lambda}_{\min}^{\epsilon})^{5/2}} \\ &= \frac{\sqrt{2} C^{2} p^{5/2}}{\sqrt{\mathbb{E}}\|\boldsymbol{Y}\|_{1}} \frac{\tilde{\lambda}_{\max}^{\star}}{(\tilde{\lambda}_{\min}^{\epsilon})^{5/2}}. \end{split}$$
(33)

The first inequality in (33) is obtained by combining the bounds in (30) and (32). The second inequality in (33) is using the fact that

$$\|\boldsymbol{y}\|_{1} \leq \mathbb{E} \|\boldsymbol{Y}\|_{1} + \frac{1}{2}\mathbb{E} \|\boldsymbol{Y}\|_{1} \leq 2\mathbb{E} \|\boldsymbol{Y}\|_{1}$$

in the event $y \in \mathcal{E}_1$. Moreover, a union bound shows that

$$\begin{split} \mathbb{P}(\Delta \not\in \mathcal{R}) &\leq \mathbb{P}(\mathcal{E}_1^c) + \mathbb{P}(\mathcal{E}_2^c) \\ &\leq \exp\left(-2\,p\right) + \frac{5}{\mathbb{E}\,\|\boldsymbol{Y}\|_1} + \frac{8\,[D_g]^+}{(\mathbb{E}\,\|\boldsymbol{Y}\|_1)^2} \\ &\leq \exp\left(-2\,p\right) + \frac{5+8\,C_0}{\mathbb{E}\,\|\boldsymbol{Y}\|_1}, \end{split}$$

where the constant C_0 and the last inequality follow from Assumption 1. Hence,

$$\mathbb{P}\left(\|\Delta\|_2 \leq \frac{\sqrt{2} C^2 p^{5/2}}{\sqrt{\mathbb{E}}\|\boldsymbol{Y}\|_1} \frac{\widetilde{\lambda}_{\max}^{\star}}{(\widetilde{\lambda}_{\min}^{\epsilon})^{5/2}}\right) \geq 1 - \exp\left(-2p\right) - \frac{5 + 8C_0}{\mathbb{E}}\|\boldsymbol{Y}\|_1.$$

Taken together, we have shown under the assumptions of Theorem 1 that there exists $N_0 \geq 3$ such that, for all $N \geq N_0$, the error of the multivariate normal approximation

$$\left|\mathbb{P}((I(\boldsymbol{\theta}^{\star}) \| \boldsymbol{Y} \|_{1})^{1/2} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^{\star}) - \Delta \in \mathcal{A}) - \Phi(\boldsymbol{Z} \in \mathcal{A})\right|$$

is bounded above by

$$\frac{83}{(\widetilde{\lambda}_{\min}^{\epsilon})^{3/2}}\sqrt{\frac{p^{7/2}}{\mathbb{E}\,\|\boldsymbol{Y}\|_{1}}} + \frac{4}{\mathbb{E}\,\|\boldsymbol{Y}\|_{1}} + \frac{8\,\left[D_{g}\right]^{+}}{\left(\mathbb{E}\,\|\boldsymbol{Y}\|_{1}\right)^{2}}$$

where Δ satisfies

$$\mathbb{P}\left(\|\Delta\|_{2} \leq \frac{\sqrt{2} C^{2} p^{5/2}}{\sqrt{\mathbb{E}}\|\boldsymbol{Y}\|_{1}} \frac{\widetilde{\lambda}_{\max}^{\star}}{(\widetilde{\lambda}_{\min}^{\epsilon})^{5/2}}\right) \geq 1 - \exp\left(-2 p\right) - \frac{5 + 8 C_{0}}{\mathbb{E}}\|\boldsymbol{Y}\|_{1}.$$

G.1. Auxiliary results for proof of Theorem 3

Lemma 5. Consider a separable multilayer network model following the form of equation (1) and is defined on a set of $N \ge 3$ and $K \ge 1$ layers and denote by $\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ the log-likelihood function. Then, for each $i = 1, \ldots, p$ and $\epsilon > 0$,

$$\sup_{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2 \le \epsilon} \quad \left| \frac{\partial^2 \left(\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \right)_i}{\partial \, \theta_j \, \partial \, \theta_k} \right| \le 2 \, \|\boldsymbol{y}\|_1,$$

where $(\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}))_i$ is the *i*th component of the score function $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$.

PROOF OF LEMMA 5. By Proposition 1, given the observation x of X (i.e., observation of the event X = x), Y is predictable with unique value $y \in \mathbb{Y}$ given by the formula in Proposition 1, and (x, y) is network concordant. Further, by Proposition 1

$$\ell(oldsymbol{ heta}; oldsymbol{x}, oldsymbol{y}) \ = \ \log \mathbb{P}_{oldsymbol{ heta}}(oldsymbol{X} = oldsymbol{x} \mid oldsymbol{Y} = oldsymbol{y}) + \log \, g(oldsymbol{y}),$$

where log $\mathbb{P}_{\boldsymbol{\theta}}(\boldsymbol{X} = \boldsymbol{x} | \boldsymbol{Y} = \boldsymbol{y})$ is the log-likelihood of a minimal, full, and regular exponential family. Thus, the second order derivative of $\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ with respect to the i^{th} and j^{th} components of $\boldsymbol{\theta}$ correspond to the variance (in the case i = j) or covariance (in the case of $i \neq j$) of corresponding sufficient statistic(s) of the exponential family [e.g., Proposition 3.8, p. 29, 41], with sufficient statistics given in Lemma 3. For $\{i, j\} \subseteq \{1, \ldots, p\}$,

$$\frac{\partial \left(\nabla_{\boldsymbol{\theta}} \,\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\right)_{i}}{\partial \,\theta_{j}} \quad = \quad \frac{\partial^{2} \,\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})}{\partial \,\theta_{i} \,\partial \,\theta_{j}} \quad = \quad \mathbb{C}_{\boldsymbol{\theta}}(s_{i}(\boldsymbol{X}), s_{j}(\boldsymbol{X}) \,|\, \boldsymbol{Y} = \boldsymbol{y}),$$

and when $i = j \in \{1, ..., p\},\$

$$\frac{\partial \left(\nabla_{\boldsymbol{\theta}} \,\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})\right)_{i}}{\partial \,\theta_{i}} \quad = \quad \frac{\partial^{2} \,\ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})}{\partial \,\theta_{i}^{2}} \quad = \quad \mathbb{V}_{\boldsymbol{\theta}}(s_{i}(\boldsymbol{X}) \,|\, \boldsymbol{Y} = \boldsymbol{y}).$$

As a result, for $\{i, j\} \subseteq \{1, \ldots, p\}$ and $k \in \{1, \ldots, p\}$,

$$\left| \frac{\partial^2 \left(\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \right)_i}{\partial \, \theta_j \, \partial \, \theta_k} \right| = \left| \frac{\partial \, \mathbb{C}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}), s_j(\boldsymbol{X}) \, | \, \boldsymbol{Y} = \boldsymbol{y})}{\partial \, \theta_k} \right|$$

and when $i = j \in \{1, ..., p\}$ and $k \in \{1, ..., p\}$,

$$\left| \frac{\partial^2 \left(\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \right)_i}{\partial \, \theta_i \, \partial \, \theta_k} \right| = \left| \frac{\partial \, \mathbb{V}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}) \, | \, \boldsymbol{Y} = \boldsymbol{y})}{\partial \, \theta_k} \right|$$

By Lemma 3 equation (15), conditional on $\mathbf{Y} = \mathbf{y}$, each sufficient statistic $s_i(\mathbf{X})$ $(i \in \{1, \ldots, p\})$ can be decomposed into the sum of conditionally independent statistics of each dyad $\mathbf{X}_{v,w}$, for $\{v, w\} \subseteq \mathbb{N}$. We can then write

$$\mathbb{C}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}), s_j(\boldsymbol{X}) \mid \boldsymbol{Y} = \boldsymbol{y}) = \sum_{\{v, w\} \subset \mathcal{N}} \mathbb{C}_{\boldsymbol{\theta}}(s_{i, v, w}(\boldsymbol{X}_{v, w}), s_{j, v, w}(\boldsymbol{X}_{v, w}) \mid \boldsymbol{Y} = \boldsymbol{y}),$$

noting that by conditional independence $\mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,r,t}(\boldsymbol{X}_{r,t}) | \boldsymbol{Y} = \boldsymbol{y}) = 0$ whenever $\{r, t\} \neq \{v, w\}$, and when i = j, we can write

$$\mathbb{V}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}) \,|\, \boldsymbol{Y} = \boldsymbol{y}) = \sum_{\{v,w\} \subset \mathbb{N}} \mathbb{V}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}) \,|\, \boldsymbol{Y} = \boldsymbol{y}),$$

again appealing to the conditional independence given Y of the random variables $s_{i,v,w}(X_{v,w})$ ($\{v,w\} \subset \mathbb{N}$). As a result, for $k \in \{1,\ldots,p\}$, it suffices to show that,

$$\frac{\partial \mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y})}{\partial \theta_k} \middle| \leq 2$$

and

$$\left| \frac{\partial \mathbb{V}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y})}{\partial \theta_k} \right| \leq 1.$$

Recall that the sufficient statistic $s_{i,v,w}(\mathbf{X})$ (i = 1, ..., p) is defined in Lemma 3 by

$$s_{i,v,w}(\boldsymbol{X}_{v,w}) = \prod_{t=1}^{h} X_{v,w}^{(k_t)}, \quad \{v,w\} \subset \mathcal{N},$$

for some $h \in \{1, \ldots, H\}$ and $\{k_1, \ldots, k_h\} \subseteq \{1, \ldots, K\}$. Define the set $S_{i,v,w}$ of components of the sufficient statistic vector $s_{v,w}(\mathbf{X})$ for $\{v, w\} \subset \mathbb{N}$ and $i = 1, \ldots, p$ by

$$S_{i,v,w} := \left\{ \prod_{t=1}^{h'} X_{v,w}^{(l_t)} : \{l_1, \dots, l_{h'}\} \subset \{k_1, \dots, k_h\}, \ h' < h \right\},\$$

where $h \in \{1, \ldots, H\}$ and $\{k_1, \ldots, k_h\} \subseteq \{1, \ldots, K\}$. The set $S_{i,v,w}$ is the set of components of the sufficient statistic vector $\mathbf{s}_{v,w}(\mathbf{X})$ of dyad $\{v, w\} \subset \mathbb{N}$ that have a value of 1 when $s_{i,v,w}(\mathbf{X}) = 1$. For the ease of notation, let $I_{S_{i,v,w}}$ be the set of dimension indices whose corresponding components of the sufficient statistic vector $\mathbf{s}_{v,w}(\mathbf{X})$ belong to the set $S_{i,v,w}$:

$$I_{S_{i,v,w}} \coloneqq \{j \in \{1, \dots, p\} : s_{j,v,w}(\mathbf{X}) \in S_{i,v,w}\}.$$

Define the conditional expectation of $s_{i,v,w}(X)$ given Y = y for any $i \in \{1, \ldots, p\}$ and $\{v, w\} \subset \mathbb{N}$ by

$$P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) := \mathbb{P}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}) = 1 | \boldsymbol{Y} = \boldsymbol{y}).$$

Denote by L_i the set of layer indices $\{k_1, \ldots, k_h\} \subseteq \{1, \ldots, K\}$ that define the i^{th} component $s_{i,v,w}(\mathbf{X}_{v,w})$ of the sufficient statistic vector $\mathbf{s}_{v,w}(\mathbf{X})$ for any $\{v, w\} \subset \mathcal{N}, j \in \{1, \ldots, p\}$, and some $h \in \{1, \ldots, H\}$. We then define

$$oldsymbol{X}_{v,w}^{(L_i)} \coloneqq \left\{ X_{v,w}^{(k_1)}, \ldots, X_{v,w}^{(k_h)}
ight\}, \quad oldsymbol{X}_{v,w}^{(-L_i)} \coloneqq oldsymbol{X}_{v,w} \setminus oldsymbol{X}_{v,w}^{(L_i)},$$

and the corresponding sample space

$$\mathbb{X}_{v,w}^{(L_i)} \coloneqq \{0,1\}^h, \qquad \mathbb{X}_{v,w}^{(-L_i)} \coloneqq \{0,1\}^{H-h},$$

for some $h \in \{1, \ldots, H\}$. Then we can write

$$P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) = \mathbb{P}_{\boldsymbol{\theta}} \left(\prod_{l \in L_i} X_{v,w}^{(l)} = 1 | \boldsymbol{Y} = \boldsymbol{y} \right)$$
$$= \frac{\sum_{\mathbb{X}_{v,w}^{(-L_i)}} \exp\left(\sum_{j \in I_{S_{i,v,w}}} \theta_j + \sum_{j \in I_{S_{i,v,w}}^c} \theta_j s_{j,v,w}(\boldsymbol{x}) \right)}{\sum_{\mathbb{X}_{v,w}} \exp\left(\sum_{j=1}^p \theta_j s_{j,v,w}(\boldsymbol{x}) \right)}.$$

Let

$$Z(\boldsymbol{\theta}) := \sum_{\mathbb{X}_{v,w}} \exp\left(\sum_{j=1}^{p} \theta_{j} s_{j,v,w}(\boldsymbol{x})\right),$$

and take the derivative of $P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y})$ with respect to θ_k for $k = 1, \ldots, p$. We have

$$\frac{\partial P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y})}{\partial \theta_{k}} \leq \frac{\sum_{\boldsymbol{X}_{v,w}^{(-L_{i})}} \exp\left(\sum_{j \in I_{S_{i,v,w}}} \theta_{j} + \sum_{j \in I_{S_{i,v,w}}^{c}} \theta_{j} s_{j,v,w}(\boldsymbol{x})\right) \left(Z(\boldsymbol{\theta}) - \frac{\partial Z(\boldsymbol{\theta})}{\partial \theta_{k}}\right)}{Z(\boldsymbol{\theta})^{2}} \\ = \frac{\sum_{\boldsymbol{X}_{v,w}^{(-L_{i})}} \exp\left(\sum_{j \in I_{S_{i,v,w}}} \theta_{j} + \sum_{j \in I_{S_{i,v,w}}^{c}} \theta_{j} s_{j,v,w}(\boldsymbol{x})\right) \left(\sum_{\boldsymbol{X}_{v,w}} \exp\left(\sum_{j=1}^{p} \theta_{j} s_{j,v,w}(\boldsymbol{x})\right) (1 - s_{k,v,w}(\boldsymbol{x}))\right)}{Z(\boldsymbol{\theta})^{2}} \\ \leq \frac{\sum_{\boldsymbol{X}_{v,w}^{(-L_{i})}} \exp\left(\sum_{j \in I_{S_{i,v,w}}} \theta_{j} + \sum_{j \in I_{S_{i,v,w}}^{c}} \theta_{j} s_{j,v,w}(\boldsymbol{x})\right) Z(\boldsymbol{\theta})}{Z(\boldsymbol{\theta})^{2}} \\ \leq 1.$$

The first inequality is obtained because $s_{k,v,w}(\boldsymbol{x}) \leq 1$, and the last inequality is due to the fact that

$$\sum_{\mathbb{X}_{v,w}^{(-L_i)}} \exp\left(\sum_{j \in I_{S_{i,v,w}}} \theta_j + \sum_{j \in I_{S_{i,v,w}}^c} \theta_j s_{j,v,w}(\boldsymbol{x})\right)$$

is bounded above by

$$\sum_{\mathbb{X}_{v,w}} \exp\left(\sum_{j=1}^p \theta_j s_{j,v,w}(\boldsymbol{x})\right).$$

Now we turn to show the derivative of the conditional variance and covariance of the sufficient statistics of each dyad are bounded. Given $\mathbf{Y} = \mathbf{y}$, for all $\{i\} \subset \{1, \ldots, p\}, s_{i,v,w}(\mathbf{X})$ are conditionally independent across $\{v, w\} \subseteq \mathcal{N}$. Then we have

$$\mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y})$$

$$= \mathbb{E}\left[s_{i,v,w}(\boldsymbol{X}) s_{j,v,w}(\boldsymbol{X}) | \boldsymbol{Y} = \boldsymbol{y}\right] - \mathbb{E}\left[s_{i,v,w}(\boldsymbol{X}) | \boldsymbol{Y} = \boldsymbol{y}\right] \mathbb{E}\left[s_{j,v,w}(\boldsymbol{X}) | \boldsymbol{Y} = \boldsymbol{y}\right]$$

$$= \mathbb{P}_{\boldsymbol{\theta}}\left(s_{i,v,w}(\boldsymbol{X}) = 1, s_{j,v,w}(\boldsymbol{X}) = 1 | \boldsymbol{Y} = \boldsymbol{y}\right) - P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) P_{j,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y})$$

$$= \mathbb{P}_{\boldsymbol{\theta}}\left(\prod_{l \in L_{i} \cup L_{j}} \boldsymbol{X}_{v,w}^{(l)} = 1 | \boldsymbol{Y} = \boldsymbol{y}\right) - P_{i,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}) P_{j,v,w}(\boldsymbol{\theta}; \boldsymbol{X}, \boldsymbol{y}).$$

Using the inequality derived in (34) and suppressing the notation of $\{v, w\}$ and (\mathbf{X}, \mathbf{y}) in $P_{i,v,w}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y})$, the derivative of the covariance with respect to θ_k , $k = 1, \ldots, p$ is given by

$$\left| \frac{\partial \mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y})}{\partial \theta_{k}} \right|$$

$$= \frac{\partial \mathbb{P}_{\boldsymbol{\theta}} \left(\prod_{l \in L_{i} \cup L_{j}} \boldsymbol{X}_{v,w}^{(l)} = 1 | \boldsymbol{Y} = \boldsymbol{y} \right)}{\partial \theta_{k}} - \frac{\partial P_{i}(\boldsymbol{\theta})}{\partial \theta_{k}} P_{j}(\boldsymbol{\theta}) - P_{i}(\boldsymbol{\theta}) \frac{\partial P_{j}(\boldsymbol{\theta})}{\partial \theta_{k}}$$

$$\leq 2.$$

Using the same inequality and notation in (34), the derivative of the variance of a Bernoulli random variable $s_{i,v,w}(\mathbf{X})$ is given by

$$\left| \frac{\partial \mathbb{V}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y})}{\partial \theta_k} \right| = \left| (1 - 2P_i(\boldsymbol{\theta})) \frac{\partial P_i(\boldsymbol{\theta})}{\partial \theta_k} \right| \leq 1.$$

Finally, for $\{i, j\} \subseteq \{1, \dots, p\}$ and $k \in \{1, \dots, p\}$, we obtain

$$\left| \frac{\partial^2 \left(\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \right)_i}{\partial \, \theta_j \, \partial \, \theta_k} \right| = \left| \frac{\partial \, \mathbb{C}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}), s_j(\boldsymbol{X}) \, | \, \boldsymbol{Y} = \boldsymbol{y})}{\partial \, \theta_k} \right|$$

$$\leq \sum_{\{v, w\} \subset \mathcal{N}} \left| \frac{\partial \, \mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,v,w}(\boldsymbol{X}_{v,w}) \, | \, \boldsymbol{Y} = \boldsymbol{y})}{\partial \, \theta_k} \right|$$

$$\leq 2 \, \| \boldsymbol{y} \|_1$$

due to the fact that $\mathbb{C}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{v,w}), s_{j,v,w}(\boldsymbol{X}_{v,w}) | \boldsymbol{Y} = \boldsymbol{y}) = 0$ when $Y_{v,w} = 0$ for $\{v,w\} \subset \mathbb{N}$. Similarly, $\mathbb{V}_{\boldsymbol{\theta}}(s_{i,v,w}(\boldsymbol{X}_{i,v,w}) | \boldsymbol{Y} = \boldsymbol{y}) = 0$ when $Y_{v,w} = 0$ for $\{v,w\} \subset \mathbb{N}$, and when $i = j \in \{1, \ldots, p\}$ and $k \in \{1, \ldots, p\}$, we have

$$\left| \frac{\partial^2 \left(\nabla_{\boldsymbol{\theta}} \, \ell(\boldsymbol{\theta}; \boldsymbol{x}, \boldsymbol{y}) \right)_i}{\partial \, \theta_i \, \partial \, \theta_k} \right| = \left| \frac{\partial \, \mathbb{V}_{\boldsymbol{\theta}}(s_i(\boldsymbol{X}) \, | \, \boldsymbol{Y} = \boldsymbol{y})}{\partial \, \theta_k} \right|$$
$$\leq \| \boldsymbol{y} \|_{1}.$$

TABLE 4
<i>P</i> -values of the Zhou-Shao's test for multivariate normality of $\tilde{\theta}$ for 6 model-generating
parameters $(\theta_1^{\star}, \theta_2^{\star}, \theta_3^{\star}, \theta_4^{\star}, \theta_5^{\star}, \theta_6^{\star})$ estimated from 250 network samples at size 1000 on
four basis network structures. All p-values are larger than .1.

Basis network model	$oldsymbol{ heta}_1^\star$	$oldsymbol{ heta}_2^{\star}$	θ_3^{\star}	$ heta_4^{\star}$	$oldsymbol{ heta}_5^{\star}$	θ_6^{\star}
Dense Bernoulli	.138	.473	.053	.699	.587	.983
Sparse Bernoulli	.554	.132	.232	.634	.904	.373
SBM	.650	.891	.982	.975	.871	.674
LSM	.859	.831	.500	.227	.613	.409

Appendix H: Additional simulation results

Additional simulation results that enhance those contained in Section 5 are provided in this section.

H.1. Normal approximation with different basis networks

The multivariate normality of $\tilde{\theta}$ is tested by Zhou-Shao's multivariate normal test [46], and the p-values are provided in tabel 4. Q-Q plots of $\tilde{\theta}$ estimated from 6 different model-generating parameters with a dense Bernoulli basis network, a sparse Bernoulli basis network, a stochastic block model (SBM) generated basis network, and a latent space model (LSM) generated basis network are shown in Figure 6, 7, 8 and 9, respectively.



Fig 6: Q-Q plots and p-values of six components of $\tilde{\theta}$ estimated from 250 multilayer network samples at size 1000 on the dense Bernoulli basis network for 6 model-generating parameters on each row.

30



Fig 7: Q-Q plots and p-values of six components of $\tilde{\theta}$ estimated from 250 multilayer network samples at size 1000 on the sparse Bernoulli basis network for 6 model-generating parameters on each row.



Fig 8: Q-Q plots and p-values of six components of $\tilde{\theta}$ estimated from 250 multilayer network samples at size 1000 on the SBM generated basis network for 6 model-generating parameters on each row.

False discovery rates of four procedures for detecting non-zero effects of six model-generating parameters $(\theta_1^{\star}, \theta_2^{\star}, \theta_3^{\star}, \theta_4^{\star}, \theta_5^{\star}, \theta_6^{\star})$ estimated from 250 multilayer network samples at size 1000 on the sparse Bernoulli basis network. All FDRs are smaller than 0.05.

Procedure	$oldsymbol{ heta}_1^\star$	$oldsymbol{ heta}_2^{\star}$	$ heta_3^{\star}$	$oldsymbol{ heta}_4^\star$	$oldsymbol{ heta}_5^{\star}$	$ heta_6^{\star}$
Bonferroni	.002	.003	.003	.003	.003	.011
Benjamini-Hochberg	.020	.011	.022	.022	.014	.017
Hochberg's	.009	.008	.012	.010	.010	.014
Holm's	.007	.008	.011	.009	.006	.014



Fig 9: Q-Q plots and p-values of six components of $\tilde{\theta}$ estimated from 250 multilayer network samples at size 1000 on the LSM generated basis network for 6 model-generating parameters on each row.

H.2. Additional results on the false discovery rate

The false discovery rate (FDR) of the multiple testing correction procedures of Bonferroni, Benjamini-Hochberg, Hochberg, and Holm to detect non-zero components of θ^* at a family-wise significance level of $\alpha = 0.05$ with a sparse Bernoulli basis network, an SBM generated basis network and an LSM generated basis network are provided in Table 5, 6 and 7, respectively (recall that components $\theta^*_{1,3}$ and θ^*_3 of θ^* are set to 0). The receiver operating characteristic (ROC) curves for $\tilde{\theta}$ of 6 selected model-generating parameters on four basis network structures are provided in each of the subplot of Figure 10.

32

TABLE 6

False discovery rates of four procedures for detecting non-zero effects of six model-generating parameters $(\theta_1^{\star}, \theta_2^{\star}, \theta_3^{\star}, \theta_4^{\star}, \theta_5^{\star}, \theta_6^{\star})$ estimated from 250 multilayer network samples at size 1000 on the SBM generated basis network. All FDRs are smaller than 0.05.

-						
Procedure	$oldsymbol{ heta}_1^\star$	$oldsymbol{ heta}_2^\star$	$ heta_3^{\star}$	$oldsymbol{ heta}_4^\star$	$oldsymbol{ heta}_5^\star$	θ_6^{\star}
Bonferroni	.002	.002	.003	.001	.001	.004
Benjamini-Hochberg	.022	.013	.014	.015	.015	.018
Hochberg's	.009	.014	.01	.008	.011	.014
Holm's	.009	.013	.005	.009	.009	.011

TABLE 7False discovery rates of four procedures for detecting non-zero effects of six model-generatingparameters ($\theta_1^*, \theta_2^*, \theta_3^*, \theta_4^*, \theta_5^*, \theta_6^*$) estimated from 250 multilayer network samples at size1000 on the LSM generated basis network. All FDRs are smaller than 0.05.

Procedure	$oldsymbol{ heta}_1^\star$	$oldsymbol{ heta}_2^{\star}$	$ heta_3^{\star}$	$oldsymbol{ heta}_4^\star$	$oldsymbol{ heta}_5^{\star}$	$oldsymbol{ heta}_6^{\star}$
Bonferroni	.004	.006	.000	.005	.003	.004
Benjamini-Hochberg	.016	.013	.011	.015	.016	.017
Hochberg's	.009	.014	.009	.011	.010	.011
Holm's	.008	.014	.009	.011	.007	.010



Fig 10: ROC curves for $\hat{\theta}$ estimated from 250 multilayer network samples at size 1000 of six model-generating parameters on four different basis networks.

References

- Albert, R. and Barabási, A. L. [2002], 'Statistical mechanics of complex networks', *Reviews of Modern Physics* 74(1), 47.
- [2] Arroyo, J. A., Athreya, A., Cape, J., Chen, G., Priebe, C. E. and Vogelstein, J. T. [2021], 'Inference for multiple heterogeneous networks with a common invariant subspace', *The Journal of Machine Learning Research* 22(1), 6303–6351.
- [3] Athreya, A., Fishkind, D. E., Tang, M., Priebe, C. E., Park, Y., Vogelstein, J. T., Levin, K., Lyzinski, V., Qin, Y. and Sussman, D. L. [2018], 'Statistical inference on random dot product graphs: a survey', *Journal of Machine Learning Research* 18(226), 1–92.
- [4] Besag, J. [1974], 'Spatial interaction and the statistical analysis of lattice systems', Journal of the Royal Statistical Society, Series B 36, 192–225.
- [5] Bhamidi, S., Bresler, G. and Sly, A. [2011], 'Mixing time of exponential random graphs', *The Annals of Applied Probability* 21, 2146–2170.
- [6] Block, P. [2015], 'Reciprocity, transitivity, and the mysterious three-cycle', Social Networks 40, 163–173.
- [7] Brown, L. [1986], Fundamentals of Statistical Exponential Families: With Applications in Statistical Decision Theory, Institute of Mathematical Statistics, Hayworth, CA, USA.
- [8] Butts, C. T. [2020], 'A dynamic process interpretation of the sparse ERGM reference model', *Journal of Mathematical Sociology*.
- [9] Caimo, A. and Gollini, I. [2020], 'A multilayer exponential random graph modelling approach for weighted networks', *Computational Statistics & Data Analysis* 142, 106825.
- [10] Caron, F. and Fox, E. B. [2017], 'Sparse graphs using exchangeable random measures', Journal of the Royal Statistical Society, Series B (with discussion) 79, 1–44.
- [11] Chen, S., Liu, S. and Ma, Z. [2022], 'Global and individualized community detection in inhomogeneous multilayer networks', *The Annals of Statistics* 50(5), 2664–2693.
- [12] Crane, H. and Dempsey, W. [2018], 'Edge exchangeable models for interaction networks', Journal of the American Statistical Association 113(523), 1311–1326.
- [13] Frank, O. [1980], 'Transitivity in stochastic graphs and digraphs', Journal of Mathematical Sociology 7, 199–213.
- [14] Furi, M. and Martelli, M. [1991], 'On the mean value theorem, inequality, and inclusion', *The American Mathematical Monthly* **98**(9), 840–846.
- [15] Geyer, C. J. and Thompson, E. A. [1992], 'Constrained Monte Carlo maximum likelihood for dependent data', *Journal of the Royal Statistical Soci*ety, Series B 54, 657–699.
- [16] Hoff, P. D., Raftery, A. E. and Handcock, M. S. [2002], 'Latent space approaches to social network analysis', *Journal of the American Statistical Association* 97, 1090–1098.
- [17] Holland, P. W., Laskey, K. B. and Leinhardt, S. [1983], 'Stochastic block

models: some first steps', Social Networks 5, 109–137.

- [18] Holland, P. W. and Leinhardt, S. [1972], 'Some evidence on the transitivity of positive interpersonal sentiment', *American Journal of Sociology* 77, 1205–1209.
- [19] Huang, S., Weng, H. and Feng, Y. [2022], 'Spectral clustering via adaptive layer aggregation for multi-layer networks', *Journal of Computational and Graphical Statistics* pp. 1–15.
- [20] Hunter, D. R., Goodreau, S. M. and Handcock, M. S. [2008], 'Goodness of fit of social network models', *Journal of the American Statistical Association* 103, 248–258.
- [21] Hunter, D. R. and Handcock, M. S. [2006], 'Inference in curved exponential family models for networks', *Journal of Computational and Graphical Statistics* 15, 565–583.
- [22] Krivitsky, P. N., Handcock, M. S. and Morris, M. [2011], 'Adjusting for network size and composition effects in exponential-family random graph models', *Statistical Methodology* 8, 319–339.
- [23] Krivitsky, P. N., Handcock, M. S., Raftery, A. E. and Hoff, P. D. [2009], 'Representing degree distributions, clustering, and homophily in social networks with latent cluster random effects models', *Social Networks* **31**, 204– 213.
- [24] Krivitsky, P. N., Koehly, L. M. and Marcum, C. S. [2020], 'Exponentialfamily random graph models for multi-layer networks', *Psychometrika* 85(3), 630–659.
- [25] Lazega, E. [2001], The collegial phenomenon: The social mechanisms of cooperation among peers in a corporate law partnership, Oxford University Press.
- [26] Lei, J., Chen, K. and Lynch, B. [2020], 'Consistent community detection in multi-layer network data', *Biometrika* 107(1), 61–73.
- [27] Li, W., Xu, Y., Yang, J. and Tang, Z. [2012], Finding structural patterns in complex networks, in '2012 IEEE Fifth International Conference on Advanced Computational Intelligence', pp. 23–27.
- [28] Lusher, D., Koskinen, J. and Robins, G. [2013], Exponential Random Graph Models for Social Networks, Cambridge University Press, Cambridge, UK.
- [29] Maathuis, M., Drton, M., Lauritzen, S. and Wainwright, M. [2018], Handbook of graphical models, CRC Press.
- [30] MacDonald, P., Levina, E. and Zhu, J. [2022], 'Latent space models for multiplex networks with shared structure', *Biometrika* 109(3), 683–706.
- [31] McPherson, M., Smith-Lovin, L. and Cook, J. M. [2001], 'Birds of a feather: Homophily in social networks', Annual Review of Sociology 27, 415–444.
- [32] Ortega, J. M. and Rheinboldt, W. C. [2000], Iterative solution of nonlinear equations in several variables, SIAM.
- [33] Raič, M. [2019], 'A multivariate Berry-Esseen theorem with explicit constants', *Bernoulli* 25(4A), 2824–2853.
- [34] Ravikumar, P., Wainwright, M. J. and Lafferty, J. [2010], 'High-dimensional Ising model selection using l₁-regularized logistic regression', *The Annals* of Statistics 38, 1287–1319.

- [35] S. Chen, D. W. a. A. S. [2015], 'Selection and estimation for mixed graphical models', *Biometrika* 102, 47–64.
- [36] Schweinberger, M., Krivitsky, P. N., Butts, C. T. and Stewart, J. [2020], 'Exponential-family models of random graphs: Inference in finite, super, and infinite population scenarios', *Statistical Science* 35, 627–662.
- [37] Sosa, J. and Betancourt, B. [2022], 'A latent space model for multilayer network data', *Computational Statistics & Data Analysis* 169, 107432.
- [38] Stewart, J. R. and Schweinberger, M. [2021], 'Pseudo-likelihood-based *M*estimation of random graphs with dependent edges and parameter vectors of increasing dimension', *arXiv preprint arXiv:2012.07167*.
- [39] Stewart, J., Schweinberger, M., Bojanowski, M. and Morris, M. [2019], 'Multilevel network data facilitate statistical inference for curved ERGMs with geometrically weighted terms', *Social Networks* 59, 98–119.
- [40] Strauss, D. and Ikeda, M. [1990], 'Pseudolikelihood estimation for social networks', Journal of the American Statistical Association 85, 204–212.
- [41] Sundberg, R. [2019], Statistical modelling by exponential families, Vol. 12, Cambridge University Press.
- [42] Vershynin, R. [2018], *High-dimensional probability: An introduction with applications in data science*, Cambridge University Press, Cambridge, UK.
- [43] Wainwright, M. J. [2019], High-dimensional statistics: A non-asymptotic viewpoint, Vol. 48, Cambridge University Press.
- [44] Wainwright, M. J. and Jordan, M. I. [2008], 'Graphical models, exponential families, and variational inference', Foundations and Trends in Machine Learning 1, 1–305.
- [45] Zhao, P. and Yu, B. [2006], 'On model selection consistency of lasso', Journal of Machine Learning Research.
- [46] Zhou, M. and Shao, Y. [2014], 'A powerful test for multivariate normality', Journal of Applied Statistics 41(2), 351–363.