## SPECTRAL INFERENCE FOR LARGE STOCHASTIC BLOCKMODELS WITH NODAL COVARIATES

#### ANGELO MELE, LINGXIN HAO, JOSHUA CAPE, AND CAREY E. PRIEBE

ABSTRACT. In many applications of network analysis, it is important to distinguish between observed and unobserved factors affecting network structure. To this end, we develop spectral estimators for both unobserved blocks and the effect of covariates in stochastic blockmodels. On the theoretical side, we establish asymptotic normality of our estimators for the subsequent purpose of performing inference. On the applied side, we show that computing our estimator is much faster than standard variational expectation–maximization algorithms and scales well for large networks. The results in this paper provide a foundation to estimate the effect of observed covariates as well as unobserved latent community structure on the probability of link formation in networks.

#### 1. INTRODUCTION

1.1. **Motivation.** The analysis and modeling of network data has important applications in sociology, economics, public health, computer science, neuroscience, and marketing, among other areas. Often, both observed and unobserved factors contribute to the global structure of networks and the processes that generate them. For example, in social networks, factors including gender, race, and personality affect the likelihood that two people interact. In such data, race and gender are typically observed, while personality is typically unobserved. It is therefore crucial to develop ways to disentangle the effect of observed and unobserved variables on link formation.

The stochastic blockmodel (SBM) is a popular network model and workhorse in the literature on community detection (Abbe, 2018; Nowicki and Snijders, 2001; Airoldi et al., 2008). In a K-block SBM, each node belongs to one of K unobserved blocks (communities); conditional on the block assignments, links form independently as Bernoulli trials with probabilities that depends on the community memberships. Many existing works use SBMs to estimate or approximate unobserved block structure in networks. In contrast, applications involving SBMs that incorporate observed nodal attributes as covariates are comparatively few (Sweet, 2015; Choi et al., 2011; Roy et al., 2019). A possible reason is that estimation for stochastic blockmodels is computational burdensome, and including covariates in the specification imposes significant additional challenges to modelling, estimation, and inference. Exact maximum likelihood estimation is infeasible, therefore most estimation strategies rely

This version: August 18, 2019. First version: February 25, 2019. Contact: angelo.mele@jhu.edu; hao@jhu.edu; jrcape@umich.edu; cep@jhu.edu. We are grateful to Cong Mu and Jipeng Zhang for excellent research assistance. We thank Eric Auerbach, Stephane Bonhomme, and Eleonora Patacchini for comments and suggestions. Funding from the Institute of Data Intensive Engineering and Science (IDIES) at Johns Hopkins University is gratefully acknowledged. Joshua Cape also gratefully acknowledges support from NSF grant DMS-1902755.

on approximations based on expectation-maximization algorithms and variational methods (Airoldi et al., 2008; Daudin et al., 2008; Bickel et al., 2013; Latouche et al., 2012; Vu et al., 2013). However, these algorithms may converge slowly to the (approximate) solution and become impractical for networks with thousands of nodes.<sup>1</sup>

1.2. **Overview.** The goal of our paper is to develop an inference procedure for SBMs with observed nodal covariates that is computationally feasible and scales well to large networks. At a high level, our main strategy consists of formulating this goal as an inference problem in the context of the generalized random dot product graph (GRDPG) model (Athreya et al., 2018; Tang and Priebe, 2018; Tang et al., 2017; Rubin-Delanchy et al., 2018). In a GRDPG, each node is characterized by an unobserved latent position (vector), and each pair of nodes link with probability determined via a (possibly indefinite) inner product of the pair's latent positions; crucially, any SBM can be reformulated as a GRDPG where the latent positions are fixed within blocks. To address our computational goals, we turn to the use of spectral methods, which, in addition to enjoying theoretical guarantees (Tang et al., 2017), have been shown to be successful both in terms of feasibility and scalability in related settings (Zheng et al., 2015, 2016, 2017).

We provide several contributions to the literature on statistical network analysis. First, we present a GRDPG model framework that incorporates the effect of observed covariates on linkage probabilities. Second, we develop a spectral estimator for inference in stochastic blockmodels with covariates, adapting the spectral estimators developed for our new class of GRDPG models. Crucially, we obtain a new central limit theorem for the spectral estimator of the covariates' effect. Our estimator is asymptotically normal as long as the parameter(s) for covariate effect can be written as sufficiently well-behaved functions of the SBM block-specific probabilities. We provide explicit formulas for bias and variance properties of the estimator, and we show that the estimator is computationally fast, scaling well for large networks. Our method provides a statistical and algorithmic foundation for inference in a broad class of models for large network data, including networks that are relatively sparse in the sense that their average degree scales sub-linearly with network size.

Our exposition focuses on SBMs with a single binary (or discrete) observed covariate, though we emphasize that the theoretical and computational properties set forth in this work extend to settings involving multiple discrete covariates. Asymptotic normality continues to hold as long as the estimator for the effect of the covariate(s) can be expressed as suitably well-behaved a function of the SBM probabilities. The case involving continuous covariates is more complicated and an active area of contemporaneous research. Current progress on this front is being facilitated by recently investigated Latent Structure Models (LSM) (Athreya et al., forthcoming) and related ideas.

The development of our estimator depends crucially on several observations. First, a Kblock stochastic blockmodel with one binary covariate can be equivalently reformulated as a (different) 2K-block stochastic blockmodel. Second, as discussed previously, a stochastic blockmodel can be viewed as a generalized random dot product graph whose latent positions are fixed within blocks (Athreya et al., 2018; Tang and Priebe, 2018; Tang et al., 2017). The

<sup>&</sup>lt;sup>1</sup>Recent advances use further approximations and parallelization to improve computational efficiency (Roy et al., 2019; Vu et al., 2013). We do not pursue such extensions in this paper.

behavior of our spectral estimation method is tied to the asymptotic behavior of spectral estimators for SBM block probability matrix entries recently studied in Tang et al. (2017). Our asymptotic analysis provides explicit formulas for standard errors and establishes the existence of a bias term; however, this bias vanishes at rate proportional to the size of the network.<sup>2</sup>

The theoretical machinery used to perform inference extends methods developed for the analysis of latent positions network models (Athreya et al., 2018; Tang et al., 2017). In particular, we use Adjacency Spectral Embedding (ASE) for random graphs to embed the network in a low-dimensional space and to recover the latent positions of the nodes. Our method is motivated by the (verifiable) intuition that the adjacency matrix can be viewed as a (mild) perturbation of the probability matrix that generates the network data, and thus, that the eigenstructure of the adjacency matrix resembles that of the edge probability matrix (Tang and Priebe, 2018; Athreya et al., 2018). In particular, spectrally decomposing the adjacency matrix provides accurate information about the structure of sufficiently large networks (Tang et al., 2017).

In addition to providing statistical guarantees, one of the many advantages of our method is the speed of computation, obtained without sacrificing estimation accuracy. In our simulations (see Section 4) we compare our approach to the variational EM (VEM) algorithm (Daudin et al., 2008; Bickel et al., 2013), as implemented in the blockmodels package in R. Even for the simplest case of a stochastic blockmodel without covariates, our spectral method is faster by several orders of magnitude. For example, in a network with n = 5000nodes and K = 2 blocks, we can estimate the model in few seconds using our spectral method, while it takes almost 10 minutes to estimate the model using the variational EM algorithm. When we add a binary covariate, our estimator converges in less than 30 seconds, while in contrast it takes almost 10 hours when using a parallelized version of the VEM algorithm in blockmodels. Our methods are implemented in the R package grdpg available at https://github.com/meleangelo/grdpg.

We also apply our method to the study of Facebook friendship data using the Facebook 100 dataset initially collected and analyzed in Traud et al. (2012). These data contain the network of friendships and node information about 100 universities in the United States in the year 2005. We estimate a stochastic blockmodel for the network of Rice University, consisting of approximately 4000 nodes, using information on gender, dorm, and year of graduation of the users (see also Roy et al. (2019)). We find evidence of homophily by gender, as suggested by the positive effect of gender on the probability of linking.

Another way to use our models and methods is to correct for the endogeneity of the network in empirical models of network effects (Shalizi and McFowland, 2018; Goldsmith-Pinkham and Imbens, 2013; Boucher and Fortin, 2016; ?; ?). Currently most of these studies rely on an auxiliary model of network formation to capture unobserved heterogeneity that affects the outcome. Our model and computational method allow the researcher to perform this type of correction for large network data.

<sup>&</sup>lt;sup>2</sup>Since we have a closed-form expression for the bias term, in principle we can naively correct for it in estimation, using a plug-in estimate. In our simulations we find that the bias term is usually so small that the correction is not necessary, at least for networks with a few thousand nodes. On the other hand, the bias is demonstrably substantial in the empirical application to Facebook data in Section 4.

#### 2. Background and methodology

2.1. Stochastic blockmodels and generalized random dot product graphs. In a *K*-block *Stochastic Blockmodel (SBM)* nodes are randomly assigned to one of *K* blocks; conditional on the blocks, nodes form links independently. A *K*-block SBM is characterized by the  $K \times K$  matrix of probabilities  $\boldsymbol{\theta} \in [0, 1]^{K \times K}$ , where the entry  $\boldsymbol{\theta}_{k\ell}$  is the probability of a link occurring between nodes in blocks *k* and  $\ell$ . The random variables comprising  $\boldsymbol{\tau} = (\tau_1, \ldots, \tau_n)$  describe the assignments of each node to a block, and they are i.i.d., such that the probability that node *i* belongs to block *k* is  $P(\tau_i = k) = \pi_k$ , with  $\boldsymbol{\pi} = (\pi_1, \ldots, \pi_K)$ . Conditional on the assignment to blocks  $\boldsymbol{\tau}$ , the probability that nodes *i* and *j* have a link is  $\boldsymbol{P}_{ij} = \boldsymbol{\theta}_{\tau_i \tau_j}$ . We use the  $n \times n$  adjacency matrix  $\boldsymbol{A}$  to describe the network, conditioning on the unobserved blocks. According to the SBM the entries of the adjacency matrix are generated as

$$\boldsymbol{A}_{ij}|\tau_i,\tau_j, \stackrel{ind}{\sim} Bernoulli(\boldsymbol{P}_{ij}), \tag{1}$$

and we write  $(\mathbf{A}, \boldsymbol{\tau}) \sim SBM(\boldsymbol{\theta}, \boldsymbol{\pi})$  to denote the adjacency matrix drawn from a K-block SBM with probability matrix  $\boldsymbol{\theta}$  and block assignment probabilities  $\boldsymbol{\pi}$ .

Conditioning on  $\boldsymbol{\tau}$ , the likelihood of the SBM with K blocks is

$$P(\boldsymbol{A}|\boldsymbol{\tau}) = \prod_{i \leq j} \boldsymbol{P}_{ij}^{\boldsymbol{A}_{ij}} (1 - \boldsymbol{P}_{ij})^{1 - \boldsymbol{A}_{ij}} = \prod_{i \leq j} \boldsymbol{\theta}_{\tau_i \tau_j}^{\boldsymbol{A}_{ij}} (1 - \boldsymbol{\theta}_{\tau_i \tau_j})^{1 - \boldsymbol{A}_{ij}}.$$
 (2)

The Generalized Random Dot Product Graph (GRPDG) model is an alternative model for network formation with conditionally independent links. In a GRDPG, each node *i* is characterized by a *d*-dimensional vector (i.e., an *unobserved* latent position)  $\mathbf{X}_i = (X_{i1}, \ldots, X_{id}) \in \mathcal{X}_d \subseteq \mathbb{R}^d$ . The latent positions are i.i.d. draws from a distribution F with support  $\mathcal{X}_d$ , that is  $\mathbf{X}_1, \mathbf{X}_2, \ldots, \mathbf{X}_n \stackrel{iid}{\sim} F$ . Let  $\mathbf{X} = [\mathbf{X}_1, \ldots, \mathbf{X}_n]^T$  denote the matrix stacking all unobserved  $\mathbf{X}_i$  vectors by row.

Let  $d_1 \geq 1$  and  $d_2 \geq 0$  be integers, and define  $d = d_1 + d_2$ . Let  $I_{d_1,d_2}$  be a  $d \times d$  diagonal matrix containing 1's in  $d_1$  diagonal entries and -1 in the remaining  $d_2$  diagonal entries. For a *GRDPG with signature*  $(d_1, d_2)$ , the entries of the adjacency matrix  $A_{ij}$  are specified to be independent, after conditioning on the latent positions  $X_i$  and  $X_j$ , namely

$$A_{ij}|X_i, X_j, \stackrel{ind}{\sim} Bernoulli(P_{ij})$$
 (3)

with link probability given by

$$\boldsymbol{P}_{ij} = \boldsymbol{X}_i^T \boldsymbol{I}_{d_1, d_2} \boldsymbol{X}_j. \tag{4}$$

For this setting, we write  $(\boldsymbol{X}, \boldsymbol{A}) \sim GRDPG_{d_1, d_2}(F)$ .<sup>3</sup>

2.2. The relationship between SBMs and GRDPGs. A remarkable property of GRDPGs is that they encompass or approximate any conditionally independent network model. In particular, any SBM can be represented as a GRDPG with latent positions fixed within blocks. That is, the K blocks are represented by a fixed location, so that each  $X_i$  can only take values  $\boldsymbol{\nu} = [\boldsymbol{\nu}_1, \boldsymbol{\nu}_2, \dots, \boldsymbol{\nu}_K]$ . Two nodes *i* and *j* belong to the same block *k* if  $X_i = X_j = \boldsymbol{\nu}_k$ . The

<sup>&</sup>lt;sup>3</sup> It must be noted that the support  $\mathcal{X}_d$  of F, is a subset of  $\mathbb{R}^d$  such that  $\boldsymbol{x}^T \boldsymbol{I}_{d_1, d_2} \boldsymbol{y} \in [0, 1]$  for all  $\boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}_d$ .

random variables  $\boldsymbol{\tau}$  are such that  $\tau_1, \ldots, \tau_n \stackrel{iid}{\sim} Multinomial(1; \pi_1, \ldots, \pi_K)$ , with  $\boldsymbol{\pi} \in (0, 1)^K$ and  $\sum_{k=1}^K \pi_k = 1$ .<sup>4</sup> The GRDPG corresponding to model  $(\boldsymbol{A}, \boldsymbol{\tau}) \sim SBM(\boldsymbol{\theta}, \boldsymbol{\pi})$  can be obtained by an eigendecomposition of the matrix  $\boldsymbol{\theta} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{U}^T$  and by defining  $\boldsymbol{\nu}_1, \boldsymbol{\nu}_2, \ldots, \boldsymbol{\nu}_K$ as the rows of  $\boldsymbol{U}|\boldsymbol{\Sigma}|^{1/2}$ . The distribution F is  $F = \sum_{k=1}^K \pi_k \delta_{\boldsymbol{\nu}_k}$ , where  $\delta$  is the Dirac-delta; importantly, d is the rank of the block-probabilities matrix  $\boldsymbol{\theta}$ , and  $d_1, d_2$  are the number of positive and negative eigenvalues of matrix  $\boldsymbol{\theta}$ , respectively.

This paper utilizes the spectral methods for inference developed for GRDPGs to estimate SBMs (Athreya et al., 2018; Tang et al., 2017). The same relationship between SBMs and GRDPG holds for *known* link functions and  $\theta_{\tau_i\tau_j} = h(B_{\tau_i\tau_j})$ , where *h* is a known function that maps to [0, 1] and **B** is a  $K \times K$  matrix of real numbers.<sup>5</sup> For example, *h* could be the logistic function or the cumulative density function of the Gaussian distribution. Our stochastic blockmodel would have adjacency matrix **A** with elements

$$\boldsymbol{A}_{ij}|\tau_i,\tau_j \stackrel{ina}{\sim} Bernoulli\left(h(\boldsymbol{B}_{\tau_i\tau_j})\right).$$
(6)

The stochastic blockmodel can be extended to include the effect of observed covariates (Choi et al., 2011; Sweet, 2015; Roy et al., 2019). Such models allow researchers to disentangle the effect of observed and unobserved nodal heterogeneity on the probability of linking. In particular, in social science, such models are used to estimate to what extent the network exhibits homophily or heterophily. Let node *i* be characterized by an *r*-dimensional vector of observed covariates  $\mathbf{Z}_i = (Z_i^{(1)}, \ldots, Z_i^{(r)}) \in \mathcal{Z} \subseteq \mathbb{R}^r$  and let the stochastic blockmodel be

$$\boldsymbol{A}_{ij}|\tau_i,\tau_j,\boldsymbol{Z}_i,\boldsymbol{Z}_j \stackrel{ina}{\sim} Bernoulli\left(h(\boldsymbol{B}_{\tau_i\tau_j}+f(\boldsymbol{Z}_i,\boldsymbol{Z}_j;\boldsymbol{\beta}))\right),\tag{7}$$

where f is a known function,  $\beta$  is a vector of parameters, and where we allow  $Z_i$  to (possibly) depend on the latent blocks.

In this paper we will focus on the case of a single binary (or discrete) covariate  $\mathbf{Z}_i$ 's, and we will assume that the function f is an indicator variables that indicates whether i and j's covariates have the same value, i.e.,

$$f(\boldsymbol{Z}_i, \boldsymbol{Z}_j; \boldsymbol{\beta}) = \boldsymbol{\beta} \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}}.$$
(8)

Here  $\beta$  can be interpreted in terms of homophily. Namely, if  $\beta > 0$ , the probability of a link between *i* and *j* is higher when their observables  $\mathbf{Z}_i$  and  $\mathbf{Z}_j$  are the same, and the network displays homophily in the observable variable. Viceversa, when  $\beta < 0$ , we have heterophily. The extension to multiple discrete covariates has similar properties and will be discussed further below.

$$\boldsymbol{X}_i \sim \pi_1 \delta_{\boldsymbol{\nu}_1} + \pi_2 \delta_{\boldsymbol{\nu}_2} + \dots + \pi_K \delta_{\boldsymbol{\nu}_K}.$$
 (5)

<sup>&</sup>lt;sup>4</sup>Alternatively, we can think of a K-block stochastic blockmodel as a network where the  $X_i$ 's are drawn from a mixture of degenerate distributions with mass centered at  $\nu$ , i.e.,

<sup>&</sup>lt;sup>5</sup>If h is unknown we cannot in general expect to be able to accurately estimate the latent positions. See Tang et al. (2013).

Our goal is to develop a general spectral method of inference for the parameter  $\beta$  and for  $B_{\tau_i \tau_i}$  in the following stochastic blockmodel with a discrete nodal covariate:

$$\boldsymbol{A}_{ij}|\tau_i,\tau_j,\boldsymbol{Z}_i,\boldsymbol{Z}_j \stackrel{ind}{\sim} Bernoulli\left(\boldsymbol{P}_{ij}\right), \tag{9}$$

$$\boldsymbol{P}_{ij} = h \left( \boldsymbol{B}_{\tau_i \tau_j} + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}} \right).$$
(10)

We further wish to disentangle the effect of observed and unobserved heterogeneity on link probabilities. To achieve this, we need to extend results from previous work on GRDPGs and SBMs (Athreya et al., 2018; Tang and Priebe, 2018; Tang et al., 2017). In the following subsections we review some of the spectral methods we use in the paper, and we provide an example that highlights the core of our method.

2.3. Spectral methods and spectral embeddings. Estimation of SBMs for large networks, with or without observed covariates, is computationally challenging. The exact MLE problem is intractable because of the high-dimensional combinatorial problem of considering all possible partitions of the nodes in blocks (Bickel et al., 2013). Approximate methods are available, based on variational approximations (Daudin et al., 2008; Airoldi et al., 2008; Wainwright and Jordan, 2008); however, even these methods are too computationally burdensome for large networks.

We make use of spectral methods, which have been shown in the literature to scale well with network size. Our spectral method embeds the network in a low(er) dimensional space, thus reducing the dimensionality of the problem, while maintaining the geometric properties of the model. In particular we use the Adjacency Spectral Embedding (ASE) to estimate the latent positions of the GRDPG (Athreya et al., 2018). In this sense, our method can be considered a dimension-reduction tool that decreases the complexity of the data by reducing the dimensionality of the space. The intuition about the spectral method is that if P is a low-rank matrix, then we can think of the adjacency matrix A as a perturbation of P, that is  $A_{ij} = P_{ij} + E_{ij}$ , where  $E_{ij}$  is a matrix of independent stochastic perturbations.<sup>6</sup> If A and P are close enough, namely if E is small enough, then the leading eigenvalues and eigenvectors of A and P will be similar (Tang et al., 2017). As a consequence, the spectral decomposition of A will provide an estimate of the latent structure of the network, that is, the latent positions X.

Consider first the case without observed covariates. Let P be positive semidefinite and let h be the identity function h(u) = u. In this setting, we only have latent positions X, that are unobserved. If we were able to observe  $P = XX^T$ , estimation of X would be straightforward. Furthermore, we could use spectral embeddings for P by exploiting the fact that P is positive semidefinite of rank d and has spectral decomposition  $P = U_P S_P U_P^T$ , where  $S_P$  is a diagonal matrix containing the largest d eigenvalues (in absolute value) of Pand  $U_P$  is the matrix with the corresponding eigenvectors. This implies that a good estimate for X is  $\widehat{X} = U_P |S_P|^{1/2}$ , where  $|\cdot|$  denotes entrywise absolute values. The estimation problem arises because we only observe A, a perturbed version of P. The Adjacency Spectral Embedding of A into  $\mathbb{R}^d$  is then  $\widehat{X} = U_A |S_A|^{1/2}$  where  $S_A$  is a diagonal matrix containing

<sup>&</sup>lt;sup>6</sup>In the Bernoulli case,  $E_{ij}$  is a shifted Bernoulli variable, with values  $E_{ij} = 1 - P_{ij}$  with probability  $P_{ij}$  and  $E_{ij} = P_{ij}$  with probability  $1 - P_{ij}$ .

the largest d eigenvalues of A in absolute value and  $U_A$  is the matrix with the corresponding eigenvectors.

In our asymptotic results for the above setup, we use the fact that ASE estimates of latent positions X asymptotically achieve perfect clustering (moreover, are asymptotically normal) and can be identified up to multiplication by an orthogonal matrix (Athreya et al., 2018; Tang et al., 2017). This implies that asymptotically the blocks are recovered exactly (up to relabeling). The same logic and results hold for non-positive definite matrices P, allowing us to study more general stochastic blockmodels (Rubin-Delanchy et al., 2018).

2.4. Overview of the method in a 2-block SBM with one binary covariate. To illustrate the methodology and to develop intuition, we focus on the special case of a K = 2 block stochastic blockmodel with a single discrete covariate and with latent positions on the unit interval [0, 1], yielding  $d_1 = 1, d_2 = 0$ , and r = 1, where  $Z_i \in \{0, 1\}$  is a binary variable (e.g., male/female, white/nonwhite, rich/poor, etc.) and the function  $f(Z_i, Z_j; \beta) = \beta \mathbf{1}_{\{Z_i = Z_j\}}$  is an indicator for the equality of the covariates for *i* and *j*, weighted by parameter  $\beta$ . The main advantage of this approach is that we can illustrate the geometry of the method in a low-dimensional space. In our simple example, the matrix **B** is given by

$$\boldsymbol{B} = \begin{array}{cc} block_1 & block_2\\ block_2 \begin{pmatrix} p^2 & pq\\ pq & q^2 \end{pmatrix} \end{array}$$
(11)

where  $p, q \in [0, 1]$ . We can conveniently re-write the matrix **B** as a dot-product of vector  $\boldsymbol{\nu} = [p \ q]^T$ , with  $p, q \in [0, 1]$ , that is  $\boldsymbol{B} = \boldsymbol{\nu} \boldsymbol{\nu}^T$ , so that the SBM can be re-written as a random dot-product graph model with  $\boldsymbol{X}_i = p$  if *i* is in block 1,  $\boldsymbol{X}_i = q$  if *i* is in block 2. The probability of linking can then be written as

$$\boldsymbol{P}_{ij} = h\left(\boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{Z_i = Z_j\}}\right).$$
(12)

For ease of exposition the network blocks have the same probability, so  $(\pi_1, \pi_2) = (0.5, 0.5)$ and each community contains half males  $(Z_i = 1)$  and half females  $(Z_i = 0)$ . However, we note that our algorithm and the theoretical results are valid when we allow the blocks to be of different size, and the observed covariates to be correlated with the unobserved blocks.

The model specified via (12) corresponds to a 4-block stochastic blockmodel. Indeed, we have 2 unobserved blocks, that are split in two additional blocks by the observed binary variable. Therefore, the final result is a 4-block SBM. More generally, if there are K latent blocks and one binary covariates, we will have a  $\tilde{K} = 2K$ -block SBM.

The possible values of  $X_i^T X_j$  are  $\{p^2, pq, q^2\}$ . Therefore the 4-block model can be completely characterized by the  $4 \times 4$  matrix

$$\boldsymbol{B}_{Z} = \begin{array}{ccc} male_{1} & female_{1} & male_{2} & female_{2} \\ male_{1} & & & \\ female_{1} \\ male_{2} \\ female_{2} \end{array} \begin{pmatrix} p^{2} + \beta & p^{2} & pq + \beta & pq \\ p^{2} & p^{2} + \beta & pq & pq + \beta \\ pq + \beta & pq & q^{2} + \beta & q^{2} \\ pq & pq + \beta & q^{2} & q^{2} + \beta \end{array} \right).$$
(13)

The value  $h(\mathbf{B}_{Z,11}) = h(p^2 + \beta)$  is the probability that two males in block 1 form a link; on the other hand,  $h(\mathbf{B}_{Z,12}) = h(p^2)$  is the probability that a male and a female in block 1 form

a link;  $h(\mathbf{B}_{Z,31}) = h(pq + \beta)$  is the probability that two males, one in block 1 and one in block 2, form a link; and so on.

The above observations imply that, for this four block SBM, there exists a corresponding GRDPG with link probability matrix

$$\boldsymbol{P} = \boldsymbol{Y} \boldsymbol{I}_{d_1, d_2} \boldsymbol{Y}^T, \tag{14}$$

for some  $n \times d$  matrix of latent positions  $\boldsymbol{Y}$  with  $d_1 \geq 1$ ,  $d_2 \geq 0$ , and  $d = d_1 + d_2$ .

To estimate the parameter  $\beta$  and the latent positions p and q we use the following algorithmic approach.

1. We compute an eigendecomposition of the adjacency matrix  $\boldsymbol{A}$ , letting  $\boldsymbol{S}_A$  denote the matrix whose diagonal contains the largest  $\hat{d}$  eigenvalues of  $\boldsymbol{A}$  in absolute value and  $\boldsymbol{U}_A$  denote the matrix whose columns are corresponding unit norm eigenvectors.<sup>7</sup> The Adjacency Spectral Embedding (ASE) of  $\boldsymbol{A}$  gives an estimate of the latent positions of the 4-block model as

$$\widehat{\boldsymbol{Y}} = \boldsymbol{U}_A |\boldsymbol{S}_A|^{1/2},\tag{15}$$

where  $|\cdot|$  indicates the absolute value (entrywise).

2. We use  $\hat{Y}$  to estimate P as

$$\widehat{\boldsymbol{P}} = \widehat{\boldsymbol{Y}} \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{Y}}^T, \tag{16}$$

where  $\hat{d} := \hat{d}_1 + \hat{d}_2$  is the number of largest eigenvalues (in magnitude) of A beyond a prescribed threshold, and  $\hat{d}_1$  and  $\hat{d}_2$  are the number of these eigenvalues that are positive and negative, respectively.

- 3. We use a clustering procedure to assign each row of  $\hat{Y}$  to one of  $\tilde{K} = 4$  blocks. We use a Gaussian Mixture Modeling approach (GMM) and estimate the center of the clusters  $\hat{\mu} = [\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3, \hat{\mu}_4]$ , that is the means of the Gaussians from the GMM.
- 4. We compute an estimate for  $\boldsymbol{\theta}_Z$  as

$$\boldsymbol{ heta}_Z = \widehat{\boldsymbol{\mu}}^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}$$

$$= \begin{array}{cccc} male_{1} & female_{1} & male_{2} & female_{2} \\ male_{1} & female_{1} & \widehat{\mu}_{1}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{1} & \widehat{\mu}_{1}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{2} & \widehat{\mu}_{1}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{3} & \widehat{\mu}_{1}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{4} \\ \widehat{\mu}_{2}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{1} & \widehat{\mu}_{2}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{2} & \widehat{\mu}_{2}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{3} & \widehat{\mu}_{2}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{4} \\ \widehat{\mu}_{3}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{1} & \widehat{\mu}_{3}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{2} & \widehat{\mu}_{3}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{3} & \widehat{\mu}_{3}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{4} \\ \widehat{\mu}_{4}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{1} & \widehat{\mu}_{4}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{2} & \widehat{\mu}_{4}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{3} & \widehat{\mu}_{4}^{T} \mathbf{I}_{\hat{d}_{1},\hat{d}_{2}} \widehat{\mu}_{4} \end{array} \right).$$
(17)

By comparing the matrix  $\hat{\theta}_Z$  and the population matrix  $\theta_Z$ , we can assign each of the 4 blocks to the original 2 blocks. In fact, we know that the diagonal terms of  $\hat{\theta}_Z$  are estimates of  $h(p^2 + \beta)$  if the nodes belong to block 1, or  $h(q^2 + \beta)$  if nodes belong to block 2. This observation shows that we can group the 4 entries of the diagonal, by checking which values are close. In our case, we will group  $\hat{\mu}_1^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_1$  and  $\hat{\mu}_2^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_2$  in one block, and  $\hat{\mu}_3^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_3$  and  $\hat{\mu}_4^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_4$  in another block. Therefore, blocks 1

<sup>&</sup>lt;sup>7</sup>In principle, the optimal value for  $\hat{d} = rank(\mathbf{B}_Z)$ ; however we do not observe  $\mathbf{B}_Z$ , therefore we estimate  $\hat{d}$  by profile likelihood methods (Zhu and Ghodsi, 2006).

and 2 in the 4-block model are assigned to the original latent block 1, while blocks 3 and 4 are assigned to original latent block 2.

5. We estimate  $\hat{\beta}$  from the entries of the matrix  $\hat{B}_Z = h^{-1}(\hat{\theta}_Z)$ , where the inverse function  $h^{-1}$  is applied element-wise. For example we know that  $h^{-1}(\hat{\mu}_1^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_1)$  is an estimate of  $p^2 + \beta$  or  $q^2 + \beta$  (because of the invariance of the model to relabeling of the blocks). Without loss of generality, assume that  $h^{-1}(\hat{\mu}_1^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_1)$  is an estimate of  $p^2 + \beta$ ; therefore, the entry  $h^{-1}(\hat{\mu}_1^T I_{\hat{d}_1,\hat{d}_2} \hat{\mu}_2)$  is an estimate of  $p^2$ . Therefore, the point estimate of  $\beta$  is then

$$\widehat{\beta} = h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_1) - h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_2).$$
(18)

6. The latent positions p and q can be estimated from the matrix  $\widehat{B}_Z$ , by using the submatrix

$$\begin{array}{ll}
 female_1 & female_2 \\
 male_1 \begin{pmatrix} h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_2) & h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_4) \\
 male_2 \begin{pmatrix} h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_2) & h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_4) \\
 h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_2) & h^{-1}(\widehat{\boldsymbol{\mu}}_1^T \boldsymbol{I}_{\hat{d}_1, \hat{d}_2} \widehat{\boldsymbol{\mu}}_4) \\
\end{array} \right) = \begin{pmatrix} \widehat{p}^2 & \widehat{p}\widehat{q} \\
 \widehat{p}\widehat{q} & \widehat{q}^2 \end{pmatrix}.$$
(19)

The spectral embedding of this matrix provides estimates for the latent positions  $\hat{p}$  and  $\hat{q}$ , that are identified up to an orthogonal transformation.

In practice, we can estimate  $\beta$  from multiple entries of the matrix  $B_Z$ , for example  $\beta = B_{Z,11} - B_{Z,12} = B_{Z,33} - B_{Z,34}$ , and weight each estimate by the size of the blocks. This could improve the estimate, since some blocks are larger than others, so delivering more precise estimates. Our code implements this idea, which is more practical for empirical applications.

#### 3. Asymptotic theory

In this section, we derive a central limit theorem for the spectral estimator of  $\beta$ . For ease of exposition, we focus on the case of a single binary observed covariate and scalar  $\beta$ , though our method works for other specifications in which the effect of the observed covariates  $\beta$  can be written as a function of the stochastic blockmodel's probability matrix  $\boldsymbol{\theta}_Z$ . Extensions to multiple binary or discrete observed covariates are straightforward.

We desire to estimate a stochastic blockmodel with observed covariates, where

$$\tau_i \stackrel{iia}{\sim} Multinomial(1; \pi_1, \dots, \pi_K),$$
 (20)

$$\mathbf{Z}_i | \tau_i \stackrel{ind}{\sim} Bernoulli(b_{\tau_i}),$$
 (21)

$$\mathbf{A}_{ij}|\tau_i, \tau_j, \mathbf{Z}_i, \mathbf{Z}_j \stackrel{ind}{\sim} Bernoulli(\mathbf{P}_{ij}),$$
 (22)

$$\boldsymbol{P}_{ij} = h\left(\boldsymbol{B}_{\tau_i\tau_j} + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_i\}}\right).$$
(23)

We assume that the observed covariates are binary and can depend on the block assignment, that is  $\mathbf{Z}_i | \tau_i \stackrel{ind}{\sim} Bernoulli(b_{\tau_i})$ , where  $b_{\tau_i} = P(\mathbf{Z}_i = 1 | \tau_i)$ . Our asymptotic results are easily extended to the case of discrete observed covariates with three or more possible outcomes.

As explained above in the simple example, our strategy consists of rewriting the SBM as a GRDPG. First, notice that the matrix  $\boldsymbol{B}$  can be written as  $\boldsymbol{B}_{\tau_i \tau_j} = \boldsymbol{X}_i^T \boldsymbol{X}_j$ , where  $\boldsymbol{X}_i$  is a  $d \times 1$  vector of latent positions that has K possible values  $\boldsymbol{\nu}_1, \ldots, \boldsymbol{\nu}_K$ . In practice,

 $\boldsymbol{\nu} = (\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_K)$  are the centers of the K blocks  $\boldsymbol{X}$ , such that *i* and *j* belong to *unobserved* block *k* when  $\boldsymbol{X}_i = \boldsymbol{X}_j = \boldsymbol{\nu}_k$ . Let  $\boldsymbol{\tau}$  be the function that assigns nodes to unobserved blocks; then  $\tau_i = k$  if  $\boldsymbol{X}_i = \boldsymbol{\nu}_k$ . We can thus rewrite the stochastic blockmodel above as a random dot product graph with observed covariates as follows:

$$\mathbf{X}_i \stackrel{\textit{ind}}{\sim} \pi_1 \delta_{\nu_1} + \pi_2 \delta_{\nu_2} + \dots + \pi_K \delta_{\nu_K}, \tag{24}$$

$$Z_i | X_i \stackrel{ind}{\sim} Bernoulli(b_{\tau_i}),$$
 (25)

$$A_{ij}|X_i, X_j, Z_i, Z_j \stackrel{ind}{\sim} Bernoulli(P_{ij}),$$
 (26)

$$\boldsymbol{P}_{ij} = h\left(\boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}}\right).$$
(27)

We first notice that both models are stochastic blockmodels with  $\tilde{K} = 2K$  blocks, because the indicator variable  $\mathbf{1}_{\{\mathbf{Z}_i = \mathbf{Z}_j\}}$  splits each unobserved block in two blocks. The probabilities of belonging to a block k for this  $\tilde{K}$ -block SBM are denoted as  $\boldsymbol{\eta} = (\eta_1, \ldots, \eta_{\tilde{K}}) =$  $(\pi_1 \cdot b_1, \pi_1 \cdot (1 - b_1), \pi_2 \cdot b_2, \pi_2 \cdot (1 - b_2), \ldots, \pi_K \cdot b_K, \pi_K \cdot (1 - b_K))$ ; and the functions that assign nodes to blocks are  $\boldsymbol{\xi} = (\xi_1, \ldots, \xi_n)$ , such that  $\xi_i = 1$  if  $\tau_i = 1$  and  $\mathbf{Z}_i = 0$ ;  $\xi_i = 2$  if  $\tau_i = 1$ and  $\mathbf{Z}_i = 1$ ;  $\xi_i = 3$  if  $\tau_i = 2$  and  $\mathbf{Z}_i = 0$ ;  $\xi_i = 4$  if  $\tau_i = 2$  and  $\mathbf{Z}_i = 1$ ; and so on.

So we have a stochastic blockmodel  $(\mathbf{A}, \boldsymbol{\xi}, \mathbf{Z}) \sim SBM(\boldsymbol{\theta}_Z, \boldsymbol{\eta})$  with  $\widetilde{K} \times \widetilde{K}$  matrix of probabilities  $\boldsymbol{\theta}_Z$ 

$$\theta_{Z} = \begin{matrix} \tau = 1; Z = 0 & \tau = 1; Z = 1 & \tau = 2; Z = 0 & \tau = 2; Z = 1 & \cdots & \tau = K; Z = 0 & \tau = K; Z = 1 \\ \tau = 1; Z = 1 & \begin{pmatrix} h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{1} \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{2} + \beta \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{2} + \beta \right) & \cdots & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{K} + \beta \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{K} + \beta \right) \\ h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{1} \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{2} + \beta \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{2} + \beta \right) & \cdots & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{K} \right) & h \left( \boldsymbol{\nu}_{1}^{T} \boldsymbol{\nu}_{K} + \beta \right) \\ h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{2} + \beta \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{2} \right) & \cdots & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{K} + \beta \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{K} \right) \\ h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{1} \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{2} \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{2} + \beta \right) & \cdots & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{K} + \beta \right) & h \left( \boldsymbol{\nu}_{2}^{T} \boldsymbol{\nu}_{K} + \beta \right) \\ \vdots & \vdots & \vdots & \ddots & \\ h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{1} \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{2} + \beta \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{K} + \beta \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{K} \right) \\ h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{1} \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{1} + \beta \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{2} + \beta \right) & \cdots & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{K} + \beta \right) & h \left( \boldsymbol{\nu}_{K}^{T} \boldsymbol{\nu}_{K} \right) \\ \end{pmatrix} \right)$$

$$(28)$$

The stochastic blockmodel characterized by matrix  $\boldsymbol{\theta}_Z$  can be re-formulated as a GRDPG. Indeed, consider the eigendecomposition of matrix  $\boldsymbol{\theta}_Z = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{U}^T$ , and define  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_{\widetilde{K}}]$  as the rows of  $\boldsymbol{U}|\boldsymbol{\Sigma}|^{1/2}$ ; then let  $F = \sum_{k=1}^{\widetilde{K}} \eta_k \delta_{\boldsymbol{\mu}_k}$ , where  $\delta$  is the Dirac-delta; and  $d_1$  and  $d_2$  are the number of positive and negative eigenvalues of  $\boldsymbol{\theta}_Z$ , respectively. So the Generalized Random Dot Product Graph model  $(\boldsymbol{Y}, \boldsymbol{A}) \sim GRDPG_{d_1,d_2}(F)$  corresponding to our stochastic blockmodel  $(\boldsymbol{A}, \boldsymbol{\xi}, \boldsymbol{Z}) \sim SBM(\boldsymbol{\theta}_Z, \boldsymbol{\eta})$  is given by

$$\boldsymbol{Y}_{i} \stackrel{iid}{\sim} \eta_{1} \delta_{\boldsymbol{\mu}_{1}} + \dots + \eta_{\tilde{K}} \delta_{\boldsymbol{\mu}_{\tilde{K}}}$$

$$\tag{29}$$

$$\boldsymbol{A}_{ij}|\boldsymbol{Y}_i, \boldsymbol{Y}_j \stackrel{ind}{\sim} Bernoulli(\boldsymbol{Y}_i^T \boldsymbol{I}_{d_1, d_2} \boldsymbol{Y}_j)$$
(30)

where  $d_1 + d_2 = \tilde{d} = rank(\boldsymbol{\theta}_Z)$  and  $\boldsymbol{Y}$  is the  $n \times \tilde{d}$  vector of latent positions with centers  $\boldsymbol{\mu}$ .

We can now extend asymptotic results for estimation of RDPGs in Athreya et al. (2018); Tang et al. (2017) to estimate block assignments and the effect of the covariates (see Rubin-Delanchy et al. (2018) for the corresponding generalization to GRDPGs).

3.1. Main theoretical result. Because the functions  $\boldsymbol{\tau}$  that describe the assignments to blocks are unknown, the  $\widetilde{K}$  SBM model assignment functions  $\boldsymbol{\xi}$  are also unknown. Applying the Adjacency Spectral Embedding procedure, we recover an estimate  $\hat{\boldsymbol{\xi}}$ .

We prove asymptotic normality for the parameter  $\beta$ , exploiting the fact that  $\beta$  can be written as a function of the SBM probabilities, that is

$$\beta = h^{-1} \left( \boldsymbol{\theta}_{Z,11} \right) - h^{-1} \left( \boldsymbol{\theta}_{Z,12} \right) = h^{-1} \left( \boldsymbol{\nu}_1^T \boldsymbol{\nu}_1 + \beta \right) - h^{-1} \left( \boldsymbol{\nu}_1^T \boldsymbol{\nu}_1 \right).$$
(31)

If the blocks were known at the onset, we could use the estimator  $\hat{\beta} = h^{-1}(\hat{\theta}_{Z,11}) - h^{-1}(\hat{\theta}_{Z,12})$ . However, all that we have access to is the estimate  $\hat{\xi}$ , so it is crucial that this estimate be consistent. For RDPGs this is indeed the case, as one can prove that the latent blocks are recovered up to an orthogonal transformation matrix in the large n limit (Lemma 4 in Tang et al. (2017)). Therefore we can recover the parameter  $\beta$  up to relabeling of the blocks. This is summarized in the following theorem.

### THEOREM 1. Central limit theorem for $\beta$

Let  $\boldsymbol{\tau}$  be unknown and K known. Let  $\hat{\boldsymbol{\tau}} : [n] \to [K]$  be the function that assigns nodes to clusters, estimated using GMM or K-means clustering on the rows of  $\hat{\boldsymbol{Y}} = \hat{\boldsymbol{U}} |\hat{\boldsymbol{S}}|^{1/2}$ . Let function g be defined as the inverse of h, that is  $g(\cdot) = h^{-1}(\cdot)$ , with first derivative  $g'(\cdot)$ . Let  $g'(\boldsymbol{\nu}_1^T \boldsymbol{\nu}_1 + \beta) \neq 0$  and  $g'(\boldsymbol{\nu}_1^T \boldsymbol{\nu}_2) \neq 0$ . Then there exists a sequence of permutations  $\phi \equiv \phi_n$  on [K] such that the estimator  $\hat{\beta} = h^{-1}(\hat{\boldsymbol{\theta}}_{Z,\phi(1)\phi(1)}) - h^{-1}(\hat{\boldsymbol{\theta}}_{Z,\phi(1)\phi(2)})$  is asymptotically normal, that is

$$n\left(\widehat{\beta} - \beta - \frac{\widehat{\psi}_{\beta}}{n}\right) \xrightarrow{d} N(0, \widehat{\sigma}_{\beta}^2)$$
(32)

as  $n \to \infty$ . The values  $\widehat{\psi}_{\beta}$  and  $\widehat{\sigma}_{\beta}^2$  are computed in the proof.

Proof. See Appendix.

This implies that we can recover the parameter  $\beta$ , and our estimator is asymptotically unbiased. In Appendix we provide the expression for the bias term.

3.2. **Sparsity.** The previous theoretical result implicitly assumes a dense network. However, many social and economic networks of interest in applications display some degree of sparsity. This is an empirical regularity that social scientists have observed in many settings, as most people do not form many links. Economists rationalize sparsity with the fact that people have constraints on time to spend with their friends.

To take this feature of the data into account, while still allowing estimation in massive networks, we multiply the probability  $P_{ij}$  by a scalar  $\rho_n$  that governs the sparsity of the network, that is, the probability of a link between nodes *i* and *j* becomes

$$\boldsymbol{P}_{ij} = \rho_n h \left( \boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}} \right).$$
(33)

Our previous result in Theorem 1 applies to dense networks; that is when  $\rho_n \to c$  where  $c \in (0, 1]$  is a constant. For simplicity and without loss of generality, in Theorem 1 we have assumed c = 1.

In this section we consider formally the case of  $\rho_n \to 0$  as  $n \to \infty$ . We have to limit the rate of convergence for  $\rho_n$ , because the network could become too sparse, not allowing estimation. We will describe this regime a *semi-sparse*, because we will allow  $\rho_n \to 0$  but  $n\rho_n = \omega(\sqrt{n})$ , that is the average degree of the network grows sub-linearly in n.<sup>8</sup> The intuition for this restriction is that too much sparsity makes links "too rare" and therefore spectral estimation and inference are impeded by having too few observations.

#### THEOREM 2. Central limit theorem for sparse networks

Let model (23) include a sparsity coefficient  $\rho_n$ 

$$\boldsymbol{P}_{ij} = \rho_n h \left( \boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}} \right)$$
(34)

such that  $\rho_n \to 0$  and  $n\rho_n = \omega(\sqrt{n})$  as  $n \to \infty$ . Let  $\hat{\tau}$  be assignment of each node to a block, estimated using ASE and GMM (or K-means) clustering. Then there exists a sequence of permutations  $\phi \equiv \phi_n$  on [K] such that the estimator  $\hat{\beta} = h^{-1}(\hat{\theta}_{Z,\phi(1)\phi(1)}) - h^{-1}(\hat{\theta}_{Z,\phi(1)\phi(2)})$  is asymptotically normal, that is

$$n\rho_n^{1/2}\left(\hat{\beta} - \beta - \frac{\ddot{\psi}_\beta}{n\rho_n}\right) \stackrel{d}{\longrightarrow} N\left(0, \ddot{\sigma}_\beta^2\right) \tag{35}$$

where  $\ddot{\psi}_{\beta}$  and  $\ddot{\sigma}_{\beta}^2$  are computed in the proof.

Proof. See Appendix.

Theorem 2 says that as long as the network is not too sparse, the estimator of  $\beta$  will be asymptotically normal. In practical estimation exercises, we can use a rule-of-thumb procedure and check whether the number of links of the network is proportional to  $\sqrt{n}$ , the square root of the network size.

3.3. Multiple observed covariates. The asymptotic results hold for discrete observed covariates and more general models, as long as the effect of the observed covariates on the probability of linking can be written as a function of the block probabilities. Let the observed variables  $\mathbf{Z}_i = [\mathbf{Z}_i^{(1)}, \mathbf{Z}_i^{(2)}]$  include two covariates. For simplicity, we consider the case of binary variables, and we assume  $\mathbf{Z}_i^{(1)} \stackrel{ind}{\sim} Bernoulli(b_{\tau_i}^{(1)})$  and  $\mathbf{Z}_i^{(2)} \stackrel{ind}{\sim} Bernoulli(b_{\tau_i}^{(2)})$ . The results still hold for discrete variables. The model is

$$\boldsymbol{A}_{ij}|\boldsymbol{X}_i, \boldsymbol{X}_j, \boldsymbol{Z}_i, \boldsymbol{Z}_j \stackrel{\textit{ind}}{\sim} Bernoulli(\boldsymbol{P}_{ij}), \qquad (36)$$

$$\boldsymbol{P}_{ij} = h\left(\boldsymbol{X}_{i}^{T}\boldsymbol{X}_{j} + \beta_{1}\boldsymbol{1}_{\{\boldsymbol{Z}_{i}^{(1)}=\boldsymbol{Z}_{j}^{(1)}\}} + \beta_{2}\boldsymbol{1}_{\{\boldsymbol{Z}_{i}^{(2)}=\boldsymbol{Z}_{j}^{(2)}\}}\right).$$
(37)

This stochastic blockmodel has  $\tilde{K} = 4K$  blocks,  $(\boldsymbol{A}, \boldsymbol{\xi}, \boldsymbol{Z}) \sim SBM(\boldsymbol{\theta}_Z, \boldsymbol{\eta})$  with  $\tilde{K} \times \tilde{K}$  matrix of probabilities  $\boldsymbol{\theta}_Z$  given by

<sup>&</sup>lt;sup>8</sup>The notation  $n\rho_n = \omega(\sqrt{n})$  means that for any real constant a > 0 there exists an  $n_0 \ge 1$  such that  $\rho_n > a/\sqrt{n} \ge 0$  for every integer  $n \ge n_0$ .

$$\boldsymbol{\theta}_{Z} = \begin{bmatrix} \boldsymbol{W}_{11} & \boldsymbol{W}_{12} & \dots & \boldsymbol{W}_{1\tilde{K}} \\ \boldsymbol{W}_{21} & \boldsymbol{W}_{22} & \dots & \boldsymbol{W}_{2\tilde{K}} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{W}_{\tilde{K}1} & \boldsymbol{W}_{\tilde{K}2} & \dots & \boldsymbol{W}_{\tilde{K}\tilde{K}} \end{bmatrix}$$
(38)

where each matrix  $\boldsymbol{W}_{k\ell}$  is given by

The intuition is the same as the model with one covariate. The blocks can be inferred by clustering the diagonal elements of matrix  $\theta_Z$ , and the parameters  $\beta_1$  and  $\beta_2$  are functions of the  $\theta_Z$  entries, namely

$$\beta_1 = h^{-1} \left( \boldsymbol{\theta}_{Z,11} \right) - h^{-1} \left( \boldsymbol{\theta}_{Z,12} \right); \quad \beta_2 = h^{-1} \left( \boldsymbol{\theta}_{Z,11} \right) - h^{-1} \left( \boldsymbol{\theta}_{Z,13} \right).$$
(40)

As such, the main characterization of the central limit theorem holds in this case with minimal modifications.

### 4. Simulation and empirical results

4.1. Comparison with Variational EM. We compare our spectral methods to a standard algorithm used in the literature, the variational EM algorithm, as implemented in the R package blockmodels. Our methods are implemented in the package grdpg, available on Github at https://github.com/meleangelo/grdpg. We also note that the variational EM algorithm uses parallelization to increase computational efficiency, while our method is implemented without any parallelization, for networks with thousands of nodes. 4.1.1. Example 1 (No covariates). In our first example, we do not include any covariates and we assume h is the identity function, so that the probabilities are  $P_{ij} = X_i^T X_j$ . We simulate networks with n = 2000, 5000, 10000, and 20000 nodes, with latent space dimension d = 1. In Table 1 we report the results for K = 2, with block centers [p, q] = [0.1, 0.7], and matrix of probabilities

$$\boldsymbol{\theta} = \begin{bmatrix} 0.01 & 0.07\\ 0.07 & 0.49 \end{bmatrix}. \tag{41}$$

For simplicity, we assume that blocks are equally likely, that is  $(\pi_1, \pi_2) = (0.5, .0.5)$ . To evaluate the performance of the algorithms, we compare clustering accuracy and computational time. The assignment of nodes to the correct block is summarized by the Adjusted Rand Index (ARI), and the computational time is given by the CPU time in seconds. Our point estimates are shown in Table 1, and below we report the estimated block probabilities for n = 2000.

$$\widehat{\boldsymbol{\theta}}_{VEM} = \begin{bmatrix} 0.01002 & 0.07012\\ 0.07012 & 0.49086 \end{bmatrix}, \qquad \widehat{\boldsymbol{\theta}}_{GRDPG} = \begin{bmatrix} 0.00998 & 0.07001\\ 0.07001 & 0.49091 \end{bmatrix}.$$
(42)

The values  $\hat{p}$ ,  $\hat{q}$  shown in in Table 1 are obtained by singular value decomposition of the estimated probability matrix (and rotation). We notice that the VEM and GRDPG estimators produce similar point estimates and very precise clustering of the nodes, as indicated by the ARI. However, our GRDPG estimator converges much faster than the VEM. For networks with n = 10000 nodes, our method provides estimates in approximately 30 seconds, while the VEM takes more than one hour to converge to the final approximation. When n = 20000 and n = 30000, our GRDPG approach converges in about 2 minutes and less than 7 minutes, respectively, while the VEM is impractical.

Here, a crucial choice is the number of dimensions for the spectral embedding. In our simulation we know that the rank of the matrix  $\boldsymbol{\theta}$  is 1, therefore this is the optimal dimension (see Athreya et al. (2018)). We choose  $\hat{d}$  by profile likelihood methods as in Zhu and Ghodsi (2006).<sup>9</sup> The clustering of the latent positions in blocks is performed using the MCLUST method implemented in the package Mclust in R (Fraley and Raftery, 1999).

In Table 2 we report results from the same model with K = 5 and latent positions  $\boldsymbol{\nu} = (0.1, 0.3, 0.5, 0.7, 0.9)$ . The results are comparable to the previous table, our estimator scales very well with the size of the network, while obtaining the same point estimates of the VEM algorithm. In this example, the difference in scaling for the two estimators is more pronounced. In particular, going from K = 2 to K = 5 blocks does not increase the computational burden too much for the GRDPG-based estimator.

4.1.2. Example 2 (logit link and binary covariate). We consider a model with a binary nodal covariate,  $\mathbf{Z}_i \sim Bernoulli(0.5)$  and link probabilities

$$\log\left(\frac{\boldsymbol{P}_{ij}}{1-\boldsymbol{P}_{ij}}\right) = \boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}}.$$
(43)

<sup>&</sup>lt;sup>9</sup>The screeplot, not shown, displays a huge step down in the (absolute) value of the eigenvalues of the adjacency matrix at the largest eigenvalue, which suggests that 1 dimension is sufficient to approximate the structure of the adjacency matrix.

In this example we use [p,q] = [-1.5,1] and  $\beta = 1.5$ , thus the matrix  $\boldsymbol{\theta}$  is

$$logit(\boldsymbol{\theta}) = \begin{bmatrix} 2.25 & -1.5\\ -1.5 & 1 \end{bmatrix}$$
(44)

while the full matrix  $\boldsymbol{\theta}_Z$  that includes the effect of covariates is

$$logit(\boldsymbol{\theta}_Z) = \begin{bmatrix} 3.75 & 2.25 & 0.00 & -1.50 \\ 2.25 & 3.75 & -1.50 & 0.00 \\ 0.00 & -1.50 & 2.50 & 1.00 \\ -1.50 & 0.00 & 1.00 & 2.50 \end{bmatrix}.$$
 (45)

We choose  $\hat{d}$  by profile likelihood (Zhu and Ghodsi, 2006). In Figure 1 we show the screeplots. In the upper-left, we display the screeplot for the adjacency matrix, which suggest to use  $\hat{d} = 4$ . We note that the fourth largest eigenvalue is negative, and the GRDPG model takes this into account. In the center-left plot, we show the screeplot of the adjacency matrix *net of the effect of the covariates*, which we use to estimate the dimension of the unobserved latent positions  $\boldsymbol{X}$ , which suggests a dimension  $\hat{d} = 1$ .

The point estimates for  $\theta_Z$  (up to a permutation of the block labels) when the network has n = 2000 nodes are respectively

$$logit(\widehat{\boldsymbol{\theta}}_{Z,VEM}) = \begin{bmatrix} 3.7443 & 2.2410 & -0.00367 & -1.5069 \\ 2.2410 & 3.7443 & -1.5069 & -0.0036 \\ -0.00367 & -1.5069 & 2.5013 & 0.9980 \\ -1.5069 & -0.0036 & 0.9980 & 2.5013 \end{bmatrix},$$
(46)

$$\widehat{\boldsymbol{B}}_{Z} = logit(\widehat{\boldsymbol{\theta}}_{Z,GRDPG}) = \begin{bmatrix} 3.7762 & 2.2336 & -0.0002 & -1.5033 \\ 2.2336 & 3.7821 & -1.5007 & -0.0042 \\ -0.0062 & -1.5007 & 2.4979 & 0.9985 \\ -1.5095 & -0.0042 & 0.9985 & 2.5045 \end{bmatrix}.$$
(47)

According to our procedure, there are several ways to obtain an estimate of  $\beta$ . From matrix (47), we group rows 1 and 2 in one block, and rows 3 and 4 in another block by clustering the diagonal entries. We can get an estimate of  $\beta$  as  $\hat{B}_{Z,11} - \hat{B}_{Z,12}$  or  $\hat{B}_{Z,22} - \hat{B}_{Z,21}$  or  $\hat{B}_{Z,33} - \hat{B}_{Z,34}$ , etc. We know by our theorem that each of these estimators is asymptotically normal. Instead of choosing which entries to use to estimate  $\beta$ , we pool all possible estimates, weighting them by the proportion of observations that are assigned to each block. For example, the estimate  $\hat{B}_{Z,11} - \hat{B}_{Z,12}$  is weighted by the proportion of links in the network that are used to estimate it.

The point estimates for  $\beta$  reported in Table 3 are  $\hat{\beta}_{GRDPG} = 1.51201$  and  $\hat{\beta}_{VEM} = 1.50335$ . The estimated latent positions are  $\hat{p} = -1.49712$  and  $\hat{q} = 1.00067$  for the VEM; and  $\hat{p} = -1.49454$  and  $\hat{q} = 0.99926$  for the GRDPG estimator. However, it takes almost 2 hours to obtain the VEM results, while it only takes 7 seconds with our estimator. The left plots in Figure 1 show the latent positions of the GRDPG (including the effect of covariates) estimated by ASE. We plot the first coordinate against each of the other three. In the second and third plot from the top, we can notice that the latent positions nicely cluster into 4 blocks, as our theory predicts. In the bottom-left plot in Figure 1 we display the estimated



FIGURE 1. Screeplots (upper left and center left), Estimated latent positions  $\hat{Y}$  (right, only 2 dimensions out of 4 per plot) and estimated latent positions  $\hat{X}$ , that is  $\hat{p}$  and  $\hat{q}$  in Example 2 (bottom left, up to orthogonal transformation) for n = 2000.

latent positions net of the covariate effect, showing how the estimated latent positions  $\overline{X}$  cluster around the true values p and q (the black vertical lines).

As explained above, a central advantage of our approach is computational speed. Indeed, in Table 3 we show that when we increase the size of the network to n = 5000, the estimated parameters are essentially the same for VEM and GRPDG. However, the GRDPG estimator take less than 30 seconds to converge; the VEM estimate takes almost 10 hours.

4.1.3. Example 3 (logit link, binary covariate and d = 2). In Table 4 we show estimates for models with latent positions  $\nu_1 = (-1.5, -1.0)$  and  $\nu_2 = (1.0, 0.5)$ . For the simulations in the first 3 rows we set  $\beta = 1.5$ . It is quite remarkable that the computational time does not increase much, with respect to the case of d = 1.

Estimator	$\mid n$	K	p	$\hat{p}$	q	$\hat{q}$	CPU Time (s)	ARI
GRDPG	2000	2	0.1	0.09993	0.7	0.70065	1.513	1
VEM	2000	2	0.1	0.10008	0.7	0.70061	39.679	1
GRDPG	5000	2	0.1	0.10004	0.7	0.69977	8.548	1
VEM	5000	2	0.1	0.10008	0.7	0.69975	593.203	1
GRDPG	10000	2	0.1	0.09994	0.7	0.69988	32.169	1
VEM	10000	2	0.1	0.09996	0.7	0.69987	4171.218	1
GRDPG	20000	2	0.1	0.09998	0.7	0.70005	128.633	1
VEM	20000	2	NA					
GRDPG	30000	2	0.1	0.09998	0.7	0.69995	386.210	1

TABLE 1. Point Estimates and CPU time for example 1 (K = 2)

TABLE 2. Point Estimates and CPU time for example 1 (K = 5)

latent	positions/	blocks
		~ ~

	latent positions/blocks										
Estimator	n	K	0.1	0.3	0.5	0.7	0.9	CPU Time (s)	ARI		
GRDPG	2000	5	0.09976	0.30122	0.49819	0.70027	0.89987	1.543	1		
VEM	2000	5	0.09994	0.30133	0.49825	0.70018	0.89973	257.713	1		
GRDPG	5000	5	0.10015	0.29952	0.49994	0.69962	0.90003	7.982	1		
VEM	5000	5	0.10021	0.29958	0.49996	0.69958	0.89999	926.330	1		
GRDPG	10000	5	0.09982	0.29975	0.49990	0.70006	0.90006	44.659	1		
VEM	10000	5	0.09985	0.29977	0.49990	0.70004	0.90004	8128.253	1		
GRDPG	20000	5	0.10000	0.30001	0.50005	0.70019	0.89999	186.073	1		
VEM	20000	5	NA								
6											

TABLE 3. Point Estimates and CPU time for example 2 (K = 2)

Estimator	n	K	p	$\hat{p}$	q	$\hat{q}$	$\beta$	$\hat{eta}$	CPU Time (s)	ARI
GRDPG	2000	2	-1.5	-1.49744	1	1.00077	no o	covariates	4.672	1
VEM	2000	2	-1.5	-1.49712	1	1.00067	no c	covariates	48.619	1
GRDPG	2000	2	-1.5	-1.49454	1	0.99926	1.5	1.51201	7.557	1
VEM	2000	2	-1.5	-1.49712	1	1.00067	1.5	1.50335	6903.673	1
GRDPG	5000	2	-1.5	-1.50029	1	1.00030	no c	covariates	17.539	1
VEM	5000	2	-1.5	-1.50019	1	1.00024	no c	covariates	537.831	1
GRDPG	5000	2	-1.5	-1.49995	1	1.00064	1.5	1.49981	27.312	1
VEM	5000	2	-1.5	-1.50019	1	1.00024	1.5	1.49955	35331.012	1
GRDPG	10000	2	-1.5	-1.49989	1	1.00029	no c	ovariates	55.428	1
GRDPG	10000	2	-1.5	1.49992	1	0.99992	1.5	1.50190	91.067	1

Estimator	n	K	$\beta$	$\hat{eta}$	Time (s)	ARI
GRDPG	2000	2	1.5	1.51760	7.335	1
GRDPG	5000	2	1.5	1.49946	28.153	1
GRDPG	10000	2	1.5	1.50257	99.128	1
GRDPG	2000	2	0.5	0.64145	7.815	0.998
GRDPG	5000	2	0.5	0.56222	26.763	1
GRDPG	10000	2	0.5	0.51617	88.562	1

TABLE 4. Point estimates, standard errors and time for Example 3.

The second group of three rows shows the results of simulations with smaller  $\beta = 0.5$ . This makes the estimation of the covariate effect more challenging. Indeed when n = 2000 the classification in blocks and the point estimate are imprecise, as indicated by the ARI. When we increase the network size to n = 5000 and n = 10000, the accuracy of the point estimates improve significantly. This example shows that our approach is extremely useful in very large networks, where VEM may become computationally impractical.

In summary, our simple examples and simulations show that our GRDPG-based estimator is quite fast and scales well to large networks. These good computational properties are obtained without sacrificing the accuracy of the estimates, as we prove that the algorithm produces the same point estimates as the variational EM in all the examples.

4.2. Application to Facebook 100 dataset. We apply our method to study the network of Facebook friendships, using the Facebook 100 dataset from Traud et al. (2012).<sup>10</sup> We follow the analysis in Roy et al. (2019), using the Rice University network data. This network consists of 4087 nodes and 7 nodal covariates: role, gender, major, minor, dorm, year, and high school. These are all discrete variables. We focus on dorm, gender, and graduation year, and we exclude the nodes that have a missing value in any of these variables as in Roy et al. (2019). In addition, we keep only nodes whose graduation year is between 2004 and 2010 and that have at least two links.

We estimate two models with one binary covariate, focusing on gender and dorm. For the Adjacency Spectral Embedding we estimate the dimension of the latent space  $\hat{d} = 12$ , using the profile likelihood method in Zhu and Ghodsi (2006). The Gaussian mixture model clustering is performed using the MCLUST implementation of Fraley and Raftery (1999) in R. We obtain  $\hat{K} = 43$  blocks. Our procedure nicely adapts to this case of odd number of clusters. Indeed once we obtain the estimated matrix  $\hat{B}_Z$ , we cluster its diagonal to recover the (unobserved) blocks.

When we estimate the model with one binary covariate (gender), we obtain an estimate of  $\hat{\beta} = 1.78096$ . This means that in this network there is homophily by gender, that is higher probability of linking when two students share the same gender.

We also estimate the effect of living in the same dorm on the probability of forming a Facebook link. We find a  $\hat{\beta} = 1.884547$ . Not surprisingly, living in the same dorm increase the probability of becoming friends on Facebook. These results are consistent with both

<sup>&</sup>lt;sup>10</sup>The entire dataset is available at https://archive.org/details/oxford-2005-facebook-matrix.

Traud et al. (2012) and Roy et al. (2019), that found dorm to be an important determinant of linking.

#### 5. CONCLUSION

We have developed a spectral estimator for large stochastic blockmodels with nodal covariates. The main theoretical contribution is an asymptotic normality result for the estimator of the covariates' effect on the probability of linking. Our work leverages the relationship between generalized random dot product graphs and stochastic blockmodels, extending existing frameworks to include observed covariates and constructing an estimator that is fast and scalable for large networks. Our theoretical results also apply to relatively sparse graphs, which is important in a host of applications in social sciences, public health, and computer science, where network data are usually sparse.

We have shown that our method delivers the same accuracy as the variational EM algorithm, while converging much faster. Our simulations and the empirical application show that this method works best in very large networks, when the variational EM becomes impractical.

We consider the present work a first step in the study of this class of models and the foundation for inference for SBMs and other latent position models for large networks with nodal covariates. While we have focused on binary and discrete covariates in this work, extensions to continuous covariates are currently being pursued via recently developed Latent Structure Models (Athreya et al., forthcoming). In future work, similar ideas can also be applied to directed networks and bipartite networks, significantly expanding the realm of GRDPG applications.

#### References

- Abbe, Emmanuel (2018), 'Community detection and stochastic block models: Recent developments', Journal of Machine Learning Research 18(177), 1–86. URL: http://jmlr.org/papers/v18/16-480.html
- Airoldi, Edoardo, David M. Blei, Stephen E. Fienberg and Eric P. Xing (2008), 'Mixed memberships stochastic blockmodels', Journal of Machine Learning Research 9, 1981– 2014.
- Athreya, Avanti, Donniell E. Fishkind, Keith Levin, Vince Lyzinski, Youngser Park, Yichen Qin, Daniel L. Sussman, Mihn Tang, Joshua T. Vogelstein and Carey E. Priebe (2018), 'Statistical inference on random dot product graphs: A survey', *Journal of Machine Learn*ing Research 18(226), 1–92.
- Athreya, Avanti, Minh Tang, Youngser Park and Carey E. Priebe (forthcoming), 'On estimation and inference in latent structure random graphs', *Statistical Science*.
- Bickel, Peter, David Choi, Xiangyu Chang and Hai Zhang (2013), 'Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels', Ann. Statist. 41(4), 1922–1943.
- Boucher, Vintent and Bernard Fortin (2016), Oxford Handbook on the Economics of Networks, Oxford University Press, chapter Some Challenges in the Empirics of the Effects of Networks.

- Choi, David, Patrick J Wolfe and Edoardo M Airoldi (2011), 'Stochastic blockmodels with a growing number of classes', *Biometrika* **99**(2), 273–284.
- Daudin, J.-J., F. Picard and S. Robin (2008), 'A mixture model for random graphs', Statistics and Computing 18(2), 173–183.
- Fraley, C. and A. E. Raftery (1999), 'Mclust: Software for model-based cluster analysis.', Journal of Classification (16), 297–306.
- Goldsmith-Pinkham, Paul and Guido W. Imbens (2013), 'Social networks and the identification of peer effects', Journal of Business and Economic Statistics 31(3), 253–264. URL: https://doi.org/10.1080/07350015.2013.801251
- Latouche, P, E Birmele and C Ambroise (2012), 'Variational bayesian inference and complexity control for stochastic block models', *Statistical Modelling* **12**(1), 93–115.
- Nimczik, Jan Sebastian (2018), Job mobility networks and endogenous labor markets. working paper.
- Nowicki, Krzysztof and Tom A. B Snijders (2001), 'Estimation and prediction for stochastic blockstructures', *Journal of the American Statistical Association* **96**(455), 1077–1087.
- Roy, Sandipan, Yves Atchade and George Michailidis (2019), 'Likelihood inference for large scale stochastic blockmodels with covariates based on a divide-and-conquer parallelizable algorithm with communication', *Journal of Computational and Graphical Statistics* 0(0), 1–22.
- Rubin-Delanchy, Patrick, Carey E Priebe, Minh Tang and Joshua Cape (2018), A statistical interpretation of spectral embedding: the generalised random dot product graph. working paper, arXiv:1709.05506.
- Shalizi, Cosma Rohilla and Edward III McFowland (2018), Estimating causal peer influence in homophilous social networks by inferring latent locations. URL: https://arxiv.org/abs/1607.06565
- Sweet, Tracy M (2015), 'Incorporating covariates into stochastic blockmodels', Journal of Educational and Behavioral Statistics **40**(6), 635–664.
- Tang, Minh and Carey E. Priebe (2018), 'Limit theorems for eigenvectors of the normalized laplacian for random graphs.', Annals of Statistics 46, 2360–2415.
- Tang, Minh, Daniel L. Sussman and Carey E. Priebe (2013), 'Universally consistent vertex classification for latent positions graphs', *The Annals of Statistics* 41(3), 1406–1430. URL: https://doi.org/10.1214/13-AOS1112
- Tang, Minh, Joshua Cape and Carey E Priebe (2017), Asymptotically efficient estimators for stochastic blockmodels: The naive mle, the rank-constrained mle, and the spectral. working paper, https://arxiv.org/abs/1710.10936.
- Traud, A. L., P. J. Mucha and Porter M. A (2012), 'Social structure of facebook networks', *Physica A: Statistical Mechanics and its Applications* **391**(16), 4165–4180.
- Vu, Duy Q, David R Hunter and Michael Schweinberger (2013), 'Model-based clustering of large networks', The annals of applied statistics 7(2), 1010–1039.
- Wainwright, M.J. and M.I. Jordan (2008), 'Graphical models, exponential families, and variational inference', *Foundations and Trends@ in Machine Learning* 1(1-2), 1–305.
- Zheng, D., R. Burns, J. Vogelstein, C. E. Priebe and A. S. Szalay (2016), An ssd-based eigensolver for spectral analysis on billion-node graph. working paper, https://arxiv.org/abs/1602.01421.

- Zheng, Da, Disa Mhembere, Randal Burns, Joshua Vogelstein, Carey E. Priebe and Alexander S. Szalay (2015), Flashgraph: Processing billion-node graphs on an array of commodity ssds, in '13th USENIX Conference on File and Storage Technologies'.
- Zheng, Da, Disa Mhembere, Vince Lyzinski, Randal Burns, Joshua Vogelstein and Carey E. Priebe (2017), 'Semi-external memory sparse matrix multiplication for billion-node graphs', *IEEE Transactions on Parallel and Distributed Systems* 28(5), 1470–1483.
- Zhu, M. and A. Ghodsi (2006), 'Automatic dimensionality selection from the scree plot via the use of profile likelihood.', Computational Statistics and Data Analysis 51, 918–930.

#### APPENDIX A. PROOFS

We first provide the general proof strategy for the simple case in which the block assignment function  $\tau$  is known. The next two theorems provide a foundation and roadmap for the proof for the more general case.

A.1. Blocks known. The simplest case is when the latent block assignments are known, the value of d and K are known, and h is the identity function. Let  $\boldsymbol{\nu} = (\boldsymbol{\nu}_1, \ldots, \boldsymbol{\nu}_K)$  be the centers of the K blocks  $\boldsymbol{X}$ , such that i and j belong to unobserved block k when  $\boldsymbol{X}_i = \boldsymbol{X}_j = \boldsymbol{\nu}_k$ . Let  $\boldsymbol{\tau}$  be the function to assign nodes to unobserved blocks, that is  $\tau_i = k$  if  $\boldsymbol{X}_i = \boldsymbol{\nu}_k$ . For this subsection we will assume that  $\boldsymbol{\tau}$  is known. Our model is

$$\boldsymbol{A}_{ij}|\boldsymbol{X}_i, \boldsymbol{X}_j, \boldsymbol{Z}_i, \boldsymbol{Z}_j, \boldsymbol{\beta} \stackrel{ina}{\sim} Bernoulli\left(\boldsymbol{X}_i^T \boldsymbol{X}_j + \boldsymbol{\beta} \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}}\right)$$
(A.1)

with  $\mathbf{Z}_i \stackrel{ind}{\sim} Bernoulli(b_{\tau_i})$ . Therefore, model (A.1) is a  $\widetilde{K} = 2K$  SBM, because we have K unobserved blocks, each split in 2 by the observed covariates. The probabilities of belonging to a block k for this  $\widetilde{K}$ -block SBM are  $\boldsymbol{\eta} = (\eta_1, \ldots, \eta_{\widetilde{K}}) = (\pi_1 \cdot b_1, \pi_1 \cdot (1-b_1), \pi_2 \cdot b_2, \pi_2 \cdot (1-b_2), \ldots, \pi_K \cdot b_K, \pi_K \cdot (1-b_K))$ . Additionally, the assignment functions are  $\boldsymbol{\xi} = (\xi_1, \ldots, \xi_n)$ , such that  $\xi_i = 1$  if  $\tau_i = 1$  and  $\mathbf{Z}_i = 0$ ;  $\xi_i = 2$  if  $\tau_i = 1$  and  $\mathbf{Z}_i = 1$ ;  $\xi_i = 3$  if  $\tau_i = 2$  and  $\mathbf{Z}_i = 0$ ;  $\xi_i = 4$  if  $\tau_i = 2$  and  $\mathbf{Z}_i = 1$ ; and so on.

Our SBM is  $(\boldsymbol{A}, \boldsymbol{\xi}, \boldsymbol{Z}) \sim SBM(\boldsymbol{\theta}_{Z}, \boldsymbol{\eta})$  with  $\widetilde{K} \times \widetilde{K}$  matrix of probabilities  $\boldsymbol{\theta}_{Z}$ 

$$\boldsymbol{\tau} = 1; Z = 0 \quad \boldsymbol{\tau} = 1; Z = 1 \quad \boldsymbol{\tau} = 2; Z = 0 \quad \boldsymbol{\tau} = 2; Z = 1 \quad \cdots \quad \boldsymbol{\tau} = K; Z = 0 \quad \boldsymbol{\tau} = K; Z = 1$$

$$\boldsymbol{\tau} = 1; Z = 0 \quad \boldsymbol{\tau} = 1; Z = 1 \quad \boldsymbol{\tau} = 1; Z = 1; Z = 1 \quad \boldsymbol{\tau} = 1; Z = 1; Z = 1 \quad \boldsymbol{\tau} = 1; Z = 1; Z$$

The following theorem establishes asymptotic normality for the estimator  $\hat{\beta} = \hat{\theta}_{Z,11} - \hat{\theta}_{Z,12}$ .

**THEOREM A.1.** Let A be an adjacency matrix from model (A.1) with h equal to the identity function h(u) = u. Let  $\tau$  be known. Then  $\hat{\beta} = \hat{\theta}_{Z,11} - \hat{\theta}_{Z,12}$  is asymptotically normal, that is

$$n\left(\hat{\beta} - \beta - \frac{\psi_{\beta}}{n}\right) \xrightarrow{d} N(0, \sigma_{\beta}^2) \tag{A.3}$$

where both  $\psi_{\beta}$  and  $\sigma_{\beta}^2$  are derived in the appendix.

*Proof.* The results in the theorem exploit the fact that  $\beta$  is a linear function of the entries of  $\theta_Z$ , whose spectral estimators also exhibit asymptotic normality (Tang et al., 2017). There is a small bias term that goes to zero with n, and we can use the result for inference.

To prove the central limit theorem using the machinery of spectral estimation of generalized random dot product graphs models, we proceed in several steps.

## STEP 1: reformulate SBM as GRDPG

Our SBM can be thought of as a GRDPG. Indeed, consider the eigendecomposition of matrix  $\boldsymbol{\theta}_Z = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{U}^T$ , and define  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_{\widetilde{K}}]$  as the rows of  $\boldsymbol{U}|\boldsymbol{\Sigma}|^{1/2}$ , then  $F = \sum_{k=1}^{\widetilde{K}} \eta_k \delta_{\boldsymbol{\mu}_k}$ , where  $\delta$  is the Dirac-delta; and  $d_1$  and  $d_2$  are the number of positive and negative eigenvalues of  $\boldsymbol{\theta}_Z$ , respectively. So the GRDPG corresponding to our SBM  $(\boldsymbol{A}, \boldsymbol{\xi}, \boldsymbol{Z}) \sim SBM(\boldsymbol{\theta}_Z, \boldsymbol{\eta})$  is given by

$$\boldsymbol{A}_{ij}|\boldsymbol{Y}_i, \boldsymbol{Y}_j \stackrel{ind}{\sim} Bernoulli(\boldsymbol{Y}_i^T \boldsymbol{I}_{d_1, d_2} \boldsymbol{Y}_j)$$
(A.4)

where  $d_1 + d_2 = \tilde{d} = rank(\boldsymbol{\theta}_Z)$  and  $\boldsymbol{Y}$  is the  $n \times \tilde{d}$  vector of latent positions with centers  $\boldsymbol{\mu}$ .

### STEP 2: spectral estimation for GRDPG

Letting  $d = rank(\theta_Z)$ , we then perform Adjacency Spectral Embedding (ASE) for A and obtain  $A = \hat{U}\hat{S}\hat{U}^T + \hat{U}_{\perp}\hat{S}_{\perp}\hat{U}_{\perp}^T$ , where  $\hat{S}$  is the diagonal matrix containing the  $\tilde{d}$  largest eigenvalues of A in absolute value, and  $\hat{U}$  is the  $n \times \tilde{d}$  matrix whose columns are the corresponding eigenvectors of A. Using a clustering procedure we can cluster the estimated latent positions  $\hat{Y} = \hat{U}|\hat{S}|^{1/2}$  (using K-means or GMM), obtaining  $\tilde{K}$  clusters and estimates of the clusters centers  $\hat{\mu}$  and cluster assignments  $\hat{\xi}$ . Notice that these estimates are consistent (Athreya et al., 2018; Tang and Priebe, 2018).

Remember that in the present setting we assume  $\tau$  is known, so our estimates for the probabilities are

$$\hat{\boldsymbol{\theta}}_{Z,k\ell} = \hat{\boldsymbol{\mu}}_k^T \boldsymbol{I}_{d_1,d_2} \hat{\boldsymbol{\mu}}_\ell \tag{A.5}$$

for any pair  $k, \ell = 1, \ldots, \widetilde{K}$ .

# **STEP 3:** estimate $\beta$ from matrix $\hat{\theta}_Z$

From the matrix  $\boldsymbol{\theta}_{Z}$  we notice that  $\beta = \boldsymbol{\theta}_{Z,11} - \boldsymbol{\theta}_{Z,12}$ , thus we can use the estimator  $\hat{\beta} = \hat{\boldsymbol{\theta}}_{Z,11} - \hat{\boldsymbol{\theta}}_{Z,12}$ . With some algebra we obtain

$$\hat{\beta} = \widehat{\theta}_{Z,11} - \widehat{\theta}_{Z,12} \tag{A.6}$$

$$= \widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11} + \boldsymbol{\theta}_{Z,11} - \boldsymbol{\theta}_{Z,12} + \boldsymbol{\theta}_{Z,12} - \widehat{\boldsymbol{\theta}}_{Z,12}$$
(A.7)

$$= (\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) + (\boldsymbol{\theta}_{Z,11} - \boldsymbol{\theta}_{Z,12}) - (\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12})$$
(A.8)

$$= (\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) + \beta - (\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12})$$
(A.9)

$$= \beta + (\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) - (\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12}).$$
(A.10)

Multiplying by n and rearranging terms we finally get

1

$$n(\hat{\beta} - \beta) = n(\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) - n(\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12}).$$
(A.11)

Hence, understanding the asymptotic behavior of  $n(\hat{\beta} - \beta)$  is equivalent to understanding the asymptotic behavior of the difference between  $n(\hat{\theta}_{Z,11} - \theta_{Z,11})$  and  $n(\hat{\theta}_{Z,12} - \theta_{Z,12})$ . It turns out that analogously to what is available for MLE and variational approximations (Bickel et al., 2013), we can prove normality of these two terms for the GRDPG (Tang et al., 2017).

The following Lemma A.1 (corresponding to Theorem 2 in Tang et al. (2017)), showing asymptotic normality of the spectral estimator for the SBM probabilities, will be used in the proof.

**LEMMA A.1.** (Theorem 2 in Tang et al. (2017)) Let  $\mathbf{A} \sim SBM(\boldsymbol{\theta}, \boldsymbol{\eta})$  be a K-block stochastic blockmodel graph on n vertices. Let  $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K$  be point masses in  $\mathbb{R}^d$  such that  $\boldsymbol{\theta}_{k\ell} = \boldsymbol{\mu}_k^T \mathbf{I}_{d_1,d_2} \boldsymbol{\mu}_\ell$  and let  $\Delta = \sum_k \eta_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$ . For  $k \in [K]$  and  $\ell \in [K]$ , let  $\psi_{k\ell}$  be

$$\psi_{k\ell} = \sum_{r=1}^{K} \eta_r \left( \boldsymbol{\theta}_{kr} (1 - \boldsymbol{\theta}_{kr}) + \boldsymbol{\theta}_{\ell r} (1 - \boldsymbol{\theta}_{\ell r}) \right) \boldsymbol{\mu}_k^T \Delta^{-1} \boldsymbol{I}_{d_1, d_2} \Delta^{-1} \boldsymbol{\mu}_\ell$$
(A.12)

$$-\sum_{r=1}^{K}\sum_{s=1}^{K}\eta_{r}\eta_{s}\boldsymbol{\theta}_{sr}(1-\boldsymbol{\theta}_{sr})\boldsymbol{\mu}_{s}^{T}\Delta^{-1}\boldsymbol{I}_{d_{1},d_{2}}\Delta^{-1}(\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}+\boldsymbol{\mu}_{k}\boldsymbol{\mu}_{\ell}^{T})\Delta^{-1}\boldsymbol{\mu}_{s}.$$
 (A.13)

Let  $\zeta_{k\ell} = \boldsymbol{\mu}_k^T \Delta^{-1} \boldsymbol{\mu}_\ell$  and define  $\sigma_{kk}^2$  for  $k \in [K]$  to be

$$\sigma_{kk}^2 = 4\boldsymbol{\theta}_{kk}(1-\boldsymbol{\theta}_{kk})\zeta_{kk}^2 + 4\sum_{r=1}^K \eta_r \boldsymbol{\theta}_{kr}(1-\boldsymbol{\theta}_{kr})\zeta_{kr}^2 \left(\frac{1}{\eta_k} - 2\zeta_{kk}\right)$$
(A.14)

+ 
$$2\sum_{r=1}^{K}\sum_{s=1}^{K}\eta_{r}\eta_{s}\boldsymbol{\theta}_{rs}(1-\boldsymbol{\theta}_{rs})\zeta_{kr}^{2}\zeta_{ks}^{2}$$
 (A.15)

and define  $\sigma_{k\ell}^2$  for  $k \in [K]$  and  $\ell \in [K]$ ,  $k \neq \ell$  to be

$$\sigma_{k\ell}^2 = \left(\boldsymbol{\theta}_{kk}(1-\boldsymbol{\theta}_{kk}) + \boldsymbol{\theta}_{\ell\ell}(1-\boldsymbol{\theta}_{\ell\ell})\right)\zeta_{kl}^2 + 2\boldsymbol{\theta}_{k\ell}(1-\boldsymbol{\theta}_{k\ell})\zeta_{kk}\zeta_{\ell\ell} \tag{A.16}$$

+ 
$$\sum_{r=1}^{K} \eta_r \boldsymbol{\theta}_{kr} (1 - \boldsymbol{\theta}_{kr}) \zeta_{\ell r}^2 \left( \frac{1}{\eta_k} - 2\zeta_{kk} \right)$$
(A.17)

+ 
$$\sum_{r=1}^{K} \eta_r \boldsymbol{\theta}_{\ell r} (1 - \boldsymbol{\theta}_{\ell r}) \zeta_{kr}^2 \left( \frac{1}{\eta_{\ell}} - 2\zeta_{\ell \ell} \right)$$
(A.18)

$$- 2\sum_{r=1}^{K} \eta_r \left( \boldsymbol{\theta}_{kr} (1 - \boldsymbol{\theta}_{kr}) + \boldsymbol{\theta}_{\ell r} (1 - \boldsymbol{\theta}_{\ell r}) \right) \zeta_{kr} \zeta_{r\ell} \zeta_{k\ell}$$
(A.19)

+ 
$$\frac{1}{2} \sum_{r=1}^{K} \sum_{s=1}^{K} \eta_r \eta_s \boldsymbol{\theta}_{rs} (1 - \boldsymbol{\theta}_{rs}) \left( \zeta_{kr} \zeta_{\ell s} + \zeta_{\ell r} \zeta_{ks} \right)^2$$
. (A.20)

Then for any  $k \in [K]$  and  $\ell \in [K]$ ,

$$n\left(\widehat{\boldsymbol{\theta}}_{k\ell} - \boldsymbol{\theta}_{k\ell} - \frac{\psi_{k\ell}}{n}\right) \stackrel{d}{\to} N(0, \sigma_{k\ell}^2)$$
(A.21)

as  $n \to \infty$ .

*Proof.* See Tang et al. (2017) for a detailed proof.

Using the result in Lemma A.1, we can see that

$$n(\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) \stackrel{d}{\to} N(\psi_{11}, \sigma_{11}^2)$$
 (A.22)

$$n(\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12}) \stackrel{d}{\to} N(\psi_{12}, \sigma_{12}^2)$$
 (A.23)

and by consequence

$$n(\hat{\beta} - \beta) \xrightarrow{d} N(\psi_{\beta}, \sigma_{\beta}^2)$$
 (A.24)

where

$$\psi_{\beta} = \psi_{11} - \psi_{12} \tag{A.25}$$

$$\sigma_{\beta}^2 = \sigma_{11}^2 + \sigma_{12}^2 - 2\sigma_{11,12} \tag{A.26}$$

and we have used notation  $\sigma_{k\ell,k'\ell'}$  to indicated the covariance terms, that is

$$\sigma_{k\ell,k'\ell'} = \mathbb{COV}(\widehat{\theta}_{Z,k\ell}, \widehat{\theta}_{Z,k'\ell'})$$
(A.27)

for any  $k, \ell, k', \ell' \in \{1, \ldots, K\}$ . The first two terms of the variance are given above in Lemma A.1. The multivariate version of the CLT can be obtained by applying the Cramer-Wold device. For the covariance term  $\sigma_{11,12}$ , we provide the calculation below.

## Computation of the covariance terms

We compute the covariance for a slightly more general model, that includes a sparsity coefficient  $\rho_n$ , that is

$$\boldsymbol{P}_{ij} = \rho_n \left( \boldsymbol{X}_i^T \boldsymbol{X}_j + \mathbf{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}} \right).$$
(A.28)

In the main text we separately discuss the cases in which  $\rho_n \to c$ , where c > 0 is a constant, or  $\rho_n \to 0$ , but  $n\rho_n = \omega(\sqrt{n})$  when  $n \to \infty$ . Let  $s_k$  be the vector in  $\mathbb{R}^n$  whose *i*-th entry is 1 if  $\xi_i = k$  and 0 otherwise, and let  $n_k$  be the number of nodes in block k, that is  $n_k = |\{i : \xi_i = k\}|$ .

We want to compute the correlation between  $\hat{\theta}_{Z,k\ell}$  and  $\hat{\theta}_{Z,k'\ell'}$ , for  $k \neq k'$  and  $\ell \neq \ell'$  in general.

To simplify notation, we will omit the Z from the subscript, so we will refer to  $\hat{\theta}_{k\ell}$  instead of  $\hat{\theta}_{Z,k\ell}$  for any  $k, \ell$ . Let **S** be the  $d \times d$  diagonal matrix containing the largest d eigenvalues of **P** in absolute value, and let **U** be the  $n \times d$  matrix whose rows are the corresponding eigenvectors of **P**. We start from equation (A.5) in the appendix of Tang et al. (2017).

$$\frac{n\rho_n^{1/2}}{n_k n_\ell} \left( \hat{\boldsymbol{\theta}}_{k\ell} - \boldsymbol{\theta}_{k\ell} \right) = \frac{n\rho_n^{-1/2}}{n_k n_\ell} \left( \boldsymbol{s}_k^T \boldsymbol{E} \boldsymbol{\Pi}_U \boldsymbol{s}_\ell + \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U^{\perp} \boldsymbol{E} \boldsymbol{\Pi}_U \boldsymbol{s}_k \right)$$
(A.29)

+ 
$$\frac{n\rho_n^{-1/2}}{n_k n_\ell} \left( \boldsymbol{s}_k^T \boldsymbol{\Pi}_U^{\perp} \boldsymbol{E}^2 \boldsymbol{P}^{\dagger} \boldsymbol{s}_\ell + \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U^{\perp} \boldsymbol{E}^2 \boldsymbol{P}^{\dagger} \boldsymbol{s}_k \right)$$
 (A.30)

+ 
$$O_p\left(n^{-1/2}\rho_n^{-1}\right)$$
 (A.31)

where  $\boldsymbol{E} = \boldsymbol{A} - \boldsymbol{P}, \, \boldsymbol{\Pi}_U = \boldsymbol{U}\boldsymbol{U}^T, \, \boldsymbol{\Pi}_U^{\perp} = \boldsymbol{I} - \boldsymbol{\Pi}_U \text{ and } \boldsymbol{P}^{\dagger} = \boldsymbol{U}\boldsymbol{S}^{-1}\boldsymbol{U}^T.$ 

Term (A.31) goes to zero when  $n\rho_n = \omega(\sqrt{n})$  as  $n \to \infty$ .

Term (A.30) is the bias term, corresponding to  $\psi_{k\ell}$  or  $\rho_n^{-1/2}\tilde{\psi}_{k\ell}$  depending on whether  $\rho_n \equiv 1 \text{ or } \rho_n \to 0, \text{ respectively.}$ 

Term (A.29) is the leading order term which converges in distribution to a normal random variable. In particular, this is the term that must be considered when deriving asymptotic covariances.

Towards this end, define  $\Upsilon_{k\ell}$  below as in equation (A.6) found in Tang et al. (2017), namely

$$\boldsymbol{\Upsilon}_{k\ell} := \frac{n\rho_n^{-1/2}}{n_k n_\ell} \left( \boldsymbol{s}_k^T \boldsymbol{E} \boldsymbol{\Pi}_U \boldsymbol{s}_\ell + \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U^{\perp} \boldsymbol{E} \boldsymbol{\Pi}_U \boldsymbol{s}_k \right)$$
(A.32)

$$= \frac{n\rho_n^{-1/2}}{n_k n_\ell} tr \boldsymbol{E} \left( \boldsymbol{\Pi}_U \boldsymbol{s}_\ell \boldsymbol{s}_k^T + \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U^\perp \right)$$
(A.33)

$$= \frac{n\rho_n^{-1/2}}{n_k n_\ell} tr \boldsymbol{E} \left( \boldsymbol{\Pi}_U \boldsymbol{s}_\ell \boldsymbol{s}_k^T + \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T - \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U \right)$$
(A.34)

$$= \frac{n\rho_n^{-1/2}}{n_k n_\ell} tr\left(\boldsymbol{A} - \boldsymbol{P}\right) \left(\boldsymbol{\Pi}_U \boldsymbol{s}_\ell \boldsymbol{s}_k^T + \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T - \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U\right)$$
(A.35)

$$= \frac{n\rho_n^{-1/2}}{n_k n_\ell} tr\left(\boldsymbol{A} - \boldsymbol{P}\right) \boldsymbol{M}$$
(A.36)

where  $\boldsymbol{M} := \boldsymbol{\Pi}_U \boldsymbol{s}_\ell \boldsymbol{s}_k^T + \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T - \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U.$ We therefore have:

$$\boldsymbol{\Upsilon}_{k\ell} = \frac{n\rho_n^{-1/2}}{n_k n_\ell} \sum_i \sum_j \left( \boldsymbol{A}_{ij} - \boldsymbol{P}_{ij} \right) \boldsymbol{M}_{ij}.$$
(A.37)

First, note that the variable  $\Upsilon_{k\ell}$  has expected value equal to zero, i.e.,  $\mathbb{E}[\Upsilon_{k\ell}] = 0$ , since  $\mathbb{E}[\mathbf{A}_{ij}] = \mathbf{P}_{ij}$  for each pair  $\{i, j\}$ . This implies that we can focus on computing covariances as

$$\mathbb{COV}[\Upsilon_{k\ell},\Upsilon_{k'\ell'}] = \mathbb{E}[\Upsilon_{k\ell}\Upsilon_{k'\ell'}].$$
(A.38)

Towards this end, let the matrix Q be defined as

$$\boldsymbol{Q} = \boldsymbol{\Pi}_{U} \boldsymbol{s}_{\ell'} \boldsymbol{s}_{k'}^{T} + \boldsymbol{\Pi}_{U} \boldsymbol{s}_{k'} \boldsymbol{s}_{\ell'}^{T} - \boldsymbol{\Pi}_{U} \boldsymbol{s}_{k'} \boldsymbol{s}_{\ell'}^{T} \boldsymbol{\Pi}_{U}, \qquad (A.39)$$

and so

$$\boldsymbol{\Upsilon}_{k'\ell'} = \frac{n\rho_n^{-1/2}}{n_{k'}n_{\ell'}} \sum_i \sum_j \left(\boldsymbol{A}_{ij} - \boldsymbol{P}_{ij}\right) \boldsymbol{Q}_{ij} \tag{A.40}$$

The product  $\Upsilon_{k\ell}\Upsilon_{k'\ell'}$  is then given by

$$\begin{split} \boldsymbol{\Upsilon}_{k\ell} \boldsymbol{\Upsilon}_{k'\ell'} &= \left( \frac{n\rho_n^{-1/2}}{n_k n_\ell} \sum_i \sum_j \left( \boldsymbol{A}_{ij} - \boldsymbol{P}_{ij} \right) \boldsymbol{M}_{ij} \right) \left( \frac{n\rho_n^{-1/2}}{n_{k'} n_{\ell'}} \sum_{i'} \sum_{j'} \left( \boldsymbol{A}_{i'j'} - \boldsymbol{P}_{i'j'} \right) \boldsymbol{Q}_{i'j'} \right) \\ &= \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_i \sum_j \sum_{i'} \sum_{j'} \left( \boldsymbol{A}_{ij} - \boldsymbol{P}_{ij} \right) \left( \boldsymbol{A}_{i'j'} - \boldsymbol{P}_{i'j'} \right) \boldsymbol{M}_{ij} \boldsymbol{Q}_{i'j'} \end{split}$$
(A.42)

First note that when  $\{i, j\} \neq \{i', j'\}$ , then the expected value of the corresponding term in the above summation is zero, since then  $(\mathbf{A}_{ij} - \mathbf{P}_{ij})$  and  $(\mathbf{A}_{i'j'} - \mathbf{P}_{i'j'})$  are independent, centered random variables. Therefore we can focus on the case when  $\{i, j\} = \{i', j'\}$ .

Further expanding  $\mathbb{E}[\Upsilon_{k\ell}\Upsilon_{k'\ell'}]$  subsequently yields

$$\mathbb{E}[\boldsymbol{\Upsilon}_{k\ell}\boldsymbol{\Upsilon}_{k'\ell'}] = \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_i \sum_j \sum_{i'} \sum_{j'} \mathbb{E}\left[ (\boldsymbol{A}_{ij} - \boldsymbol{P}_{ij}) \left( \boldsymbol{A}_{i'j'} - \boldsymbol{P}_{i'j'} \right) \right] \boldsymbol{M}_{ij} \boldsymbol{Q}_{i'j} (A.43)$$

$$= \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_i \sum_j \mathbb{E} \left[ (\boldsymbol{A}_{ij} - \boldsymbol{P}_{ij})^2 \right] \boldsymbol{M}_{ij} \boldsymbol{Q}_{ij}$$
(A.44)

$$= \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_i \sum_j \boldsymbol{P}_{ij} \left(1 - \boldsymbol{P}_{ij}\right) \boldsymbol{M}_{ij} \boldsymbol{Q}_{ij}.$$
(A.45)

We thus need to compute the entries of M and Q to obtain a computable formula. We have for M that  $M_{ij} := v_{ij}^{(1)} + v_{ij}^{(2)} + v_{ij}^{(3)}$ , where

$$v_{ij}^{(1)} := \left( \boldsymbol{\Pi}_U \boldsymbol{s}_\ell \boldsymbol{s}_k^T \right)_{ij} = n_\ell \boldsymbol{Y}_i^T (\boldsymbol{Y}^T \boldsymbol{Y})^{-1} \boldsymbol{\mu}_\ell \mathbb{1}\{\xi_j = k\}$$
(A.46)

$$v_{ij}^{(2)} := \left( \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \right)_{ij} = n_k \boldsymbol{Y}_i^T (\boldsymbol{Y}^T \boldsymbol{Y})^{-1} \boldsymbol{\mu}_k \mathbb{1}\{\xi_j = \ell\}$$
(A.47)

$$v_{ij}^{(3)} := \left( \boldsymbol{\Pi}_U \boldsymbol{s}_k \boldsymbol{s}_\ell^T \boldsymbol{\Pi}_U \right)_{ij} = n_k n_\ell \boldsymbol{Y}_i^T (\boldsymbol{Y}^T \boldsymbol{Y})^{-1} \boldsymbol{\mu}_k \boldsymbol{\mu}_\ell^T (\boldsymbol{Y}^T \boldsymbol{Y})^{-1} \boldsymbol{Y}_j$$
(A.48)

(A.49)

and analogously for  $\boldsymbol{Q}_{ij} := \varrho_{ij}^{(1)} + \varrho_{ij}^{(2)} + \varrho_{ij}^{(3)}$ . Hence,  $\boldsymbol{M}_{ij}\boldsymbol{Q}_{ij} = \sum_{1 \le \alpha,\beta \le 3} v_{ij}^{(\alpha)}\varrho_{ij}^{(\beta)}$ . In what follows, it will sometimes be useful to write the scalars  $v_{ij}^{(\alpha)}, \varrho_{ij}^{(\beta)}$  equivalently in terms of their transpose (i.e., see the right-hand sides of their definitions).

To begin, we make a preliminary observation that

$$\frac{\rho_n^{-1}}{n_k n_{k'}} \left( \sum_{i,j} \boldsymbol{P}_{ij} (1 - \boldsymbol{P}_{ij}) \boldsymbol{Y}_i \boldsymbol{Y}_i^T \mathbb{1}\{\xi_j = k\} \mathbb{1}\{\xi_j = k'\} \right)$$

$$\stackrel{a.s.}{\to} \begin{cases} \frac{1}{\eta_k} \mathbb{E} \left[ \boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_1)} (1 - \boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_1)}) \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] & \text{if } \rho_n \equiv 1 \text{ and } k = k', \\ \frac{1}{\eta_k} \mathbb{E} \left[ \boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_1)} \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] & \text{if } \rho_n \to 0 \text{ and } k = k', \\ 0 & \text{if } k \neq k'. \end{cases}$$

Consider the terms involving  $v_{ij}^{(1)} \rho_{ij}^{(1)}, v_{ij}^{(1)} \rho_{ij}^{(2)}, v_{ij}^{(2)} \rho_{ij}^{(1)}$ , and  $v_{ij}^{(2)} \rho_{ij}^{(2)}$ . Then

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\upsilon_{ij}^{(1)}\varrho_{ij}^{(1)} \\
= \frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\left(n_{\ell}n_{\ell'}\boldsymbol{\mu}_{\ell}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\left(\sum_{i,j}\boldsymbol{P}_{ij}(1-\boldsymbol{P}_{ij})\boldsymbol{Y}_{i}\boldsymbol{Y}_{i}^{T}\mathbb{1}\{\xi_{j}=k\}\mathbb{1}\{\xi_{j}=k'\}\right)(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{\ell'}\right) \\
\xrightarrow{a.s.} \begin{cases} \frac{1}{\eta_{k}}\boldsymbol{\mu}_{\ell}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})}(1-\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})})\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell'} & \text{if } \rho_{n}\equiv 1 \text{ and } k=k', \\ \frac{1}{\eta_{k}}\boldsymbol{\mu}_{\ell}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})}\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell'} & \text{if } \rho_{n}\to 0 \text{ and } k=k', \\ 0 & \text{if } k\neq k', \end{cases}$$

and similarly

$$\frac{{}^{n^{2}\rho_{n}^{-1}}}{{}^{n_{k}n_{\ell}n_{k'}n_{\ell'}}} \sum_{i,j} \boldsymbol{P}_{ij} \left(1-\boldsymbol{P}_{ij}\right) \upsilon_{ij}^{(1)} \varrho_{ij}^{(2)}$$

$$\stackrel{a.s.}{\rightarrow} \begin{cases} \frac{1}{\eta_{k}} \boldsymbol{\mu}_{\ell}^{T} \boldsymbol{\Delta}^{-1} \mathbb{E} \left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})} (1-\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})}) \boldsymbol{Y}_{1} \boldsymbol{Y}_{1}^{T}\right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_{n} \equiv 1 \text{ and } k = \ell', \\ \frac{1}{\eta_{k}} \boldsymbol{\mu}_{\ell}^{T} \boldsymbol{\Delta}^{-1} \mathbb{E} \left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})} \boldsymbol{Y}_{1} \boldsymbol{Y}_{1}^{T}\right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_{n} \to 0 \text{ and } k = \ell', \\ 0 & \text{if } k \neq \ell', \end{cases}$$

and similarly

$$\begin{split} & \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_{i,j} \boldsymbol{P}_{ij} \left( 1 - \boldsymbol{P}_{ij} \right) \upsilon_{ij}^{(2)} \varrho_{ij}^{(1)} \\ & \underset{\rightarrow}{\overset{a.s.}{\rightarrow}} \begin{cases} \frac{1}{\eta_\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)} (1 - \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)}) \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell'} & \text{if } \rho_n \equiv 1 \text{ and } \ell = k', \\ \frac{1}{\eta_\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)} \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell'} & \text{if } \rho_n \to 0 \text{ and } \ell = k', \\ 0 & \text{if } \ell \neq k'. \end{cases}$$

and similarly

$$\frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_{i,j} \boldsymbol{P}_{ij} \left( 1 - \boldsymbol{P}_{ij} \right) \upsilon_{ij}^{(2)} \varrho_{ij}^{(2)}$$

$$\stackrel{a.s.}{\rightarrow} \begin{cases} \frac{1}{\eta_\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)} (1 - \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)}) \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_n \equiv 1 \text{ and } \ell = \ell', \\ \frac{1}{\eta_\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\boldsymbol{Y}_1)} \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_n \to 0 \text{ and } \ell = \ell', \\ 0 & \text{if } \ell \neq \ell'. \end{cases}$$

Next, we consider the terms involving  $v_{ij}^{(1)} \rho_{ij}^{(3)}$ ,  $v_{ij}^{(2)} \rho_{ij}^{(3)}$ ,  $v_{ij}^{(3)} \rho_{ij}^{(1)}$ , and  $v_{ij}^{(3)} \rho_{ij}^{(2)}$ . In particular,

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\boldsymbol{v}_{ij}^{(1)}\boldsymbol{\varrho}_{ij}^{(3)} \\
= \frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\left(n_{\ell}n_{k'}n_{\ell'}\mathbb{1}\left\{\xi_{j}=k\right\}\boldsymbol{\mu}_{\ell}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{Y}_{i}\boldsymbol{Y}_{i}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{Y}_{j}\right) \\
\xrightarrow{a.s.} \begin{cases} \boldsymbol{\mu}_{\ell}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})}(1-\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})})\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k} & \text{if } \rho_{n}\equiv 1, \\ \boldsymbol{\mu}_{\ell}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k\xi(\boldsymbol{Y}_{1})}\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k} & \text{if } \rho_{n}\rightarrow 0, \end{cases}$$

and similarly

$$\begin{split} & \frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_{i,j} \boldsymbol{P}_{ij} \left( 1 - \boldsymbol{P}_{ij} \right) \boldsymbol{v}_{ij}^{(2)} \boldsymbol{\varrho}_{ij}^{(3)} \\ & \stackrel{a.s.}{\to} \begin{cases} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\mathbf{Y}_1)} (1 - \boldsymbol{\theta}_{\ell\xi(\mathbf{Y}_1)}) \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} \boldsymbol{\mu}_{\ell'}^T \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell} & \text{if } \rho_n \equiv 1, \\ \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[ \boldsymbol{\theta}_{\ell\xi(\mathbf{Y}_1)} \boldsymbol{Y}_1 \boldsymbol{Y}_1^T \right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} \boldsymbol{\mu}_{\ell'}^T \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell} & \text{if } \rho_n \equiv 0. \end{cases}$$

Along the same lines,

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\boldsymbol{v}_{ij}^{(3)}\boldsymbol{\varrho}_{ij}^{(1)} \\
= \frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\left(n_{k}n_{\ell}n_{\ell'}\mathbb{1}\left\{\xi_{j}=k'\right\}\boldsymbol{Y}_{j}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{Y}_{i}\boldsymbol{Y}_{i}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{\ell'}\right) \\
\xrightarrow{a.s.} \begin{cases} \boldsymbol{\mu}_{k'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k'\xi(\boldsymbol{Y}_{1})}(1-\boldsymbol{\theta}_{k'\xi(\boldsymbol{Y}_{1})})\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell'} & \text{if } \rho_{n}\equiv 1, \\ \boldsymbol{\mu}_{k'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{k'\xi(\boldsymbol{Y}_{1})}\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell'} & \text{if } \rho_{n}\to 0, \end{cases}$$

and similarly

$$\frac{n^2 \rho_n^{-1}}{n_k n_\ell n_{k'} n_{\ell'}} \sum_{i,j} \boldsymbol{P}_{ij} \left(1 - \boldsymbol{P}_{ij}\right) \upsilon_{ij}^{(3)} \varrho_{ij}^{(2)}$$

$$\stackrel{a.s.}{\rightarrow} \begin{cases} \boldsymbol{\mu}_{\ell'}^T \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[\boldsymbol{\theta}_{\ell'\xi(\boldsymbol{Y}_1)} (1 - \boldsymbol{\theta}_{\ell'\xi(\boldsymbol{Y}_1)}) \boldsymbol{Y}_1 \boldsymbol{Y}_1^T\right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_n \equiv 1, \\ \boldsymbol{\mu}_{\ell'}^T \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{\ell} \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \mathbb{E} \left[\boldsymbol{\theta}_{\ell'\xi(\boldsymbol{Y}_1)} \boldsymbol{Y}_1 \boldsymbol{Y}_1^T\right] \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_{k'} & \text{if } \rho_n \to 0. \end{cases}$$

Finally, consider the term involving  $v_{ij}^{(3)} \varrho_{ij}^{(3)}$ . We see that

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\upsilon_{ij}^{(3)}\varrho_{ij}^{(3)} \\
= \frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\left(n_{k}n_{\ell}n_{k'}n_{\ell'}\boldsymbol{Y}_{j}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{Y}_{i}\boldsymbol{Y}_{i}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}(\boldsymbol{Y}^{T}\boldsymbol{Y})^{-1}\boldsymbol{Y}_{j}\right)$$

In order to analyze this quantity, we now decompose the sum over the index j using the indicator variables  $\mathbb{1}{\xi_j = \alpha}$  for all SBM blocks  $\alpha \in \mathcal{A}$ . For each such term we get

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\upsilon_{ij}^{(3)}\varrho_{ij}^{(3)}\mathbb{1}\left\{\xi_{j}=\alpha\right\}$$

$$\stackrel{a.s.}{\rightarrow} \begin{cases} \eta_{\alpha}\boldsymbol{\mu}_{\alpha}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})}(1-\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})})\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\alpha} & \text{if } \rho_{n}\equiv1, \\ \eta_{\alpha}\boldsymbol{\mu}_{\alpha}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})}\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\alpha} & \text{if } \rho_{n}\neq0. \end{cases}$$

Hence, by aggregating over all  $\alpha \in \mathcal{A}$ , we obtain

$$\frac{n^{2}\rho_{n}^{-1}}{n_{k}n_{\ell}n_{k'}n_{\ell'}}\sum_{i,j}\boldsymbol{P}_{ij}\left(1-\boldsymbol{P}_{ij}\right)\upsilon_{ij}^{(3)}\varrho_{ij}^{(3)}$$

$$\stackrel{a.s.}{\rightarrow} \begin{cases} \sum_{\alpha}\eta_{\alpha}\boldsymbol{\mu}_{\alpha}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})}(1-\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})})\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\alpha} & \text{if } \rho_{n}\equiv 1, \\ \sum_{\alpha}\eta_{\alpha}\boldsymbol{\mu}_{\alpha}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\ell}\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Delta}^{-1}\mathbb{E}\left[\boldsymbol{\theta}_{\alpha\xi(\boldsymbol{Y}_{1})}\boldsymbol{Y}_{1}\boldsymbol{Y}_{1}^{T}\right]\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{k'}\boldsymbol{\mu}_{\ell'}^{T}\boldsymbol{\Delta}^{-1}\boldsymbol{\mu}_{\alpha} & \text{if } \rho_{n}\equiv 1, \end{cases}$$

Combining all of the above observations yields  $\mathbb{COV}[\Upsilon_{k\ell}\Upsilon_{k'\ell'}]$  for all possible relationships between tuples  $\{k, \ell\}, \{k', \ell'\}$  and for both regimes  $\rho_n \equiv 1, \rho_n \to 0$ .

A.2. Blocks known and general link function. If the link function is not the identity, then our general model becomes

$$\boldsymbol{A}_{ij}|\boldsymbol{X}_i, \boldsymbol{X}_j, \boldsymbol{Z}_i, \boldsymbol{Z}_j, \beta \stackrel{ind}{\sim} Bernoulli\left(h(\boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \boldsymbol{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}})\right).$$
(A.50)

A popular choice of h is the logistic specification, with  $h(u) = e^u/(1 + e^u)$  (Choi et al., 2011; Traud et al., 2012; Sweet, 2015; Nimczik, 2018). The generalization of the central limit theorem is as follows.

**THEOREM A.2.** (General nonlinear h function) Let  $\boldsymbol{A}$  be an adjacency matrix from model (A.50) and let h be a link function,  $h: \mathcal{X} \times \mathcal{X} \times \mathcal{Z} \times \mathcal{Z} \to [0,1]$ . Let  $\boldsymbol{\tau}$  be known and let function g be defined as the inverse of h, that is  $g(\cdot) = h^{-1}(\cdot)$ , with first derivative  $g'(\cdot)$ . Let  $g'(\boldsymbol{\nu}_1^T\boldsymbol{\nu}_1 + \beta) \neq 0$  and  $g'(\boldsymbol{\nu}_1^T\boldsymbol{\nu}_2) \neq 0$ . Then  $\hat{\boldsymbol{\beta}} = h^{-1}(\hat{\boldsymbol{\theta}}_{Z,11}) - h^{-1}(\hat{\boldsymbol{\theta}}_{Z,12})$  is asymptotically normal, in particular

$$n\left(\hat{\beta} - \beta - \frac{\widetilde{\psi}_{\beta}}{n}\right) \xrightarrow{d} N\left(0, \widetilde{\sigma}_{\beta}^{2}\right) \tag{A.51}$$

where  $\widetilde{\psi}_{\beta}$  and  $\widetilde{\sigma}_{\beta}^2$  are computed in appendix.

*Proof.* In this case model (23) is as follows:

$$\boldsymbol{A}_{ij}|\boldsymbol{X}_i, \boldsymbol{X}_j, \boldsymbol{Z}_i, \boldsymbol{Z}_j, \beta \stackrel{ind}{\sim} Bernoulli\left(h\left(\boldsymbol{X}_i^T \boldsymbol{X}_j + \beta \mathbf{1}_{\{\boldsymbol{Z}_i = \boldsymbol{Z}_j\}}\right)\right).$$
(A.52)

All the rest is the same, and we can write our SBM as  $(\mathbf{A}, \boldsymbol{\xi}, \mathbf{Z}) \sim SBM(\boldsymbol{\theta}_Z, \boldsymbol{\eta})$  with  $\widetilde{K} \times \widetilde{K}$  matrix of probabilities  $\boldsymbol{\theta}_Z$  given by

$$\boldsymbol{\theta}_{Z} = \begin{matrix} \tau = 1; Z = 0 & \tau = 1; Z = 1 & \tau = 2; Z = 0 & \tau = 2; Z = 1 & \cdots & \tau = K; Z = 0 & \tau = K; Z = 1 \\ \tau = 1; Z = 0 & \begin{pmatrix} h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{K} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{K} + \beta) \\ h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & \cdots & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{K} + \beta) \\ h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{2} + \beta) & \cdots & h(\boldsymbol{\nu}_{1}^{T}\boldsymbol{\nu}_{K} + \beta) \\ h(\boldsymbol{\nu}_{2}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{2}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{2}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{2}^{T}\boldsymbol{\nu}_{2} + \beta) & \cdots & h(\boldsymbol{\nu}_{2}^{T}\boldsymbol{\nu}_{K} + \beta) \\ \vdots & \vdots & \vdots & \ddots & \\ h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{2} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{2} + \beta) & \cdots & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{K} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{K} + \beta) \\ h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{1} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{2} + \beta) & \cdots & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{K} + \beta) & h(\boldsymbol{\nu}_{K}^{T}\boldsymbol{\nu}_{K} + \beta) \\ \end{pmatrix} \right) \right)$$

$$(A.53)$$

So we know that  $\beta = h^{-1}(\boldsymbol{\theta}_{Z,11}) - h^{-1}(\boldsymbol{\theta}_{Z,12})$ , and we can use the estimator

$$\hat{\boldsymbol{\beta}} = h^{-1}(\widehat{\boldsymbol{\theta}}_{Z,11}) - h^{-1}(\widehat{\boldsymbol{\theta}}_{Z,12}).$$
(A.54)

By Theorem A.1 we know that

$$n(\widehat{\boldsymbol{\theta}}_{Z,11} - \boldsymbol{\theta}_{Z,11}) \xrightarrow{d} N(\psi_{11}, \sigma_{11}^2)$$
 (A.55)

$$n(\widehat{\boldsymbol{\theta}}_{Z,12} - \boldsymbol{\theta}_{Z,12}) \stackrel{d}{\to} N(\psi_{12}, \sigma_{12}^2)$$
 (A.56)

Applying a Taylor expansion reveals

$$h^{-1}(\widehat{\theta}_{Z,11}) = h^{-1}(\theta_{Z,11}) + (h^{-1}(\theta_{Z,11}))' (\widehat{\theta}_{Z,11} - \theta_{Z,11}) + \text{smaller order terms}(A.57)$$
  
$$h^{-1}(\widehat{\theta}_{Z,12}) = h^{-1}(\theta_{Z,12}) + (h^{-1}(\theta_{Z,12}))' (\widehat{\theta}_{Z,12} - \theta_{Z,12}) + \text{smaller order terms}(A.58)$$

and this implies

$$n\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{Z,11}) - h^{-1}(\boldsymbol{\theta}_{Z,11})\right) \stackrel{d}{\to} N(\widetilde{\psi}_{11}, \widetilde{\sigma}_{11}^2), \tag{A.59}$$

$$n\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{Z,12}) - h^{-1}(\boldsymbol{\theta}_{Z,12})\right) \stackrel{d}{\to} N(\widetilde{\psi}_{12}, \widetilde{\sigma}_{12}^2), \tag{A.60}$$

where  $\widetilde{\psi}_{k\ell} = \psi_{k\ell} \left( h^{-1}(\boldsymbol{\theta}_{Z,k\ell}) \right)'$  and  $\widetilde{\sigma}_{k\ell}^2 = \sigma_{k\ell}^2 \left[ \left( h^{-1}(\boldsymbol{\theta}_{Z,k\ell}) \right)' \right]^2$  for  $k, \ell = 1, 2$ . Finally this implies that our estimator behaves in the manner

$$n(\hat{\beta} - \beta) \stackrel{d}{\to} N(\widetilde{\psi}_{\beta}, \widetilde{\sigma}_{\beta}^2),$$
 (A.61)

where the bias term is

$$\widetilde{\psi}_{\beta} = \widetilde{\psi}_{11} - \widetilde{\psi}_{12}, \tag{A.62}$$

where the variance term  $\widetilde{\sigma}_{\beta}^2$  satisfies

$$\widetilde{\sigma}_{\beta}^2 = \widetilde{\sigma}_{11}^2 + \widetilde{\sigma}_{12}^2 - 2\widetilde{\sigma}_{11,12}, \tag{A.63}$$

and where the covariance term is given by

$$\widetilde{\sigma}_{11,12} = cov(h^{-1}(\widehat{\theta}_{11}), h^{-1}(\widehat{\theta}_{12})) = \sigma_{11,12} \left[ \left( h^{-1}(\theta_{Z,11}) \right)' \right] \left[ \left( h^{-1}(\theta_{Z,12}) \right)' \right].$$
(A.64)

31

A.3. **Proof of THEOREM 1.** Let K be known. When  $\tau$  is not known, we can estimate it using Adjacency Spectral Embedding (ASE) to get an estimate  $\hat{Y} = \hat{U}|\hat{S}|^{1/2}$ ; we then cluster the rows of  $\hat{Y}$  using a Gaussian Mixture Model (GMM) or K-means clustering. This gives estimates  $\hat{\xi}$  and therefore  $\hat{\tau}$  as we can easily estimate the probability  $b_k$  of the Bernoulli covariates from the observed  $Z_i$ 's, within each estimated block.

Because of the estimation the blocks are recovered up to a permutation of the blocks' labels (as it is in any mixture model). Nonetheless, we can still obtain an asymptotic result for  $\hat{\beta}$  using the following Lemma A.2, taken from Tang et al. (2017).

**LEMMA A.2.** (Corollary 2 in Tang et al. (2017)) Let the setting and notation be as in Theorem A.1. Let K be known and let  $\hat{\boldsymbol{\xi}} : [n] \to [\tilde{K}]$  be the function that assigns nodes to clusters, estimated using GMM or K-means clustering on the rows of  $\hat{\boldsymbol{Y}} = \hat{\boldsymbol{U}}|\hat{\boldsymbol{S}}|^{1/2}$  (as in the proof of Theorem A.1). Let  $\hat{\boldsymbol{\theta}}_{Z,k\ell} = \hat{\boldsymbol{\mu}}_k^T \mathbf{I}_{d_1,d_2} \hat{\boldsymbol{\mu}}_\ell$  and let  $\hat{\boldsymbol{\Delta}} = \sum_{k=1}^K \hat{\eta}_k \hat{\boldsymbol{\mu}}_k \hat{\boldsymbol{\mu}}_\ell^T$ . For  $k \in [K]$ and  $\ell \in [K]$  define  $\hat{\psi}_{k\ell}$  as

$$\widehat{\psi}_{k\ell} = \sum_{r=1}^{K} \widehat{\xi}_r \left( \widehat{\theta}_{kr} (1 - \widehat{\theta}_{kr}) + \widehat{\theta}_{\ell r} (1 - \widehat{\theta}_{\ell r}) \right) \widehat{\mu}_k^T \widehat{\Delta}^{-1} I_{d_1, d_2} \widehat{\Delta}^{-1} \widehat{\mu}_\ell$$
(A.65)

$$-\sum_{r=1}^{K}\sum_{s=1}^{K}\widehat{\eta}_{r}\widehat{\eta}_{s}\widehat{\theta}_{sr}(1-\widehat{\theta}_{sr})\widehat{\mu}_{s}^{T}\widehat{\Delta}^{-1}\boldsymbol{I}_{d_{1},d_{2}}\widehat{\Delta}^{-1}(\widehat{\mu}_{\ell}\widehat{\mu}_{k}^{T}+\widehat{\mu}_{k}\widehat{\mu}_{\ell}^{T})\widehat{\Delta}^{-1}\widehat{\mu}_{s} \qquad (A.66)$$

Then there exists a sequence of permutations  $\phi \equiv \phi_n$  on [K] such that for any  $k \in [K]$  and  $\ell \in [K]$ ,

$$n\left(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)} - \boldsymbol{\theta}_{k\ell} - \frac{\widehat{\psi}_{k\ell}}{n}\right) \stackrel{d}{\to} N(0,\sigma_{k\ell}^2)$$
(A.67)

as  $n \to \infty$ .

*Proof.* See Tang et al. (2017) for the detailed proof.

Let us first focus on the **linear** case, in which h is the identity function. If h(u) = u, then we can estimate  $\beta$  as

$$\widehat{\beta} = \widehat{\theta}_{\phi(1),\phi(1)} - \widehat{\theta}_{\phi(1),\phi(2)}$$
(A.68)

$$= \hat{\boldsymbol{\theta}}_{\phi(1),\phi(1)} - \boldsymbol{\theta}_{11} + \boldsymbol{\theta}_{11} - \boldsymbol{\theta}_{12} + \boldsymbol{\theta}_{12} - \hat{\boldsymbol{\theta}}_{\phi(1),\phi(2)}$$
(A.69)

$$= \left(\widehat{\boldsymbol{\theta}}_{\phi(1),\phi(1)} - \boldsymbol{\theta}_{11}\right) + \beta - \left(\widehat{\boldsymbol{\theta}}_{\phi(1),\phi(2)} - \boldsymbol{\theta}_{12}\right)$$
(A.70)

and rearranging we obtain

$$n\left(\widehat{\beta}-\beta\right) = n\left(\widehat{\theta}_{\phi(1),\phi(1)} - \theta_{11}\right) - n\left(\widehat{\theta}_{\phi(1),\phi(2)} - \theta_{12}\right)$$
(A.71)

which by Lemma A.2 implies that there exists a (sequence of) permutation(s)  $\phi$  such that

$$n\left(\widehat{\beta} - \beta - \frac{\widehat{\psi}_{\beta}}{n}\right) \xrightarrow{d} N(0, \sigma_{\beta}^2) \tag{A.72}$$

where  $\widehat{\psi}_{\beta} = (\widehat{\psi}_{11} - \widehat{\psi}_{12})$  and  $\sigma_{\beta}^2 = \sigma_{11}^2 + \sigma_{12}^2 - 2cov\left(\widehat{\theta}_{\phi(1),\phi(1)}, \widehat{\theta}_{\phi(1),\phi(2)}\right)$ . For the **poplinger** link function, we use a Taylor expansion

For the **nonlinear** link function, we use a Taylor expansion

$$h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)}) = h^{-1}(\boldsymbol{\theta}_{k\ell}) + [h^{-1}(\boldsymbol{\theta}_{k\ell})]' \left(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)} - \boldsymbol{\theta}_{k\ell}\right) + \text{smaller order terms}(A.73)$$

and thus we have

$$h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)}) - h^{-1}(\boldsymbol{\theta}_{k\ell}) = [h^{-1}(\boldsymbol{\theta}_{k\ell})]' \left(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)} - \boldsymbol{\theta}_{k\ell}\right) + \text{smaller order terms}(A.74)$$

which implies that

$$n\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)}) - h^{-1}(\boldsymbol{\theta}_{k\ell})\right) \stackrel{d}{\to} N\left(\tilde{\tilde{\psi}}_{k\ell}, \tilde{\tilde{\sigma}}_{k\ell}^2\right),\tag{A.75}$$

where  $\tilde{\tilde{\psi}}_{k\ell} = [h^{-1}(\boldsymbol{\theta}_{k\ell})]' \hat{\psi}_{k\ell}$  and  $\tilde{\tilde{\sigma}}_{k\ell}^2 = ([h^{-1}(\boldsymbol{\theta}_{k\ell})]')^2 \sigma_{k\ell}^2$ . We therefore obtain the result

$$n(\widehat{\beta} - \beta) \xrightarrow{d} N(\widehat{\psi}_{\beta}, \widehat{\sigma}_{\beta}^2)$$
 (A.76)

where  $\hat{\psi}_{\beta} = \tilde{\tilde{\psi}}_{11} - \tilde{\tilde{\psi}}_{12}$  and  $\hat{\sigma}_{\beta}^2 = \tilde{\tilde{\sigma}}_{11}^2 + \tilde{\tilde{\sigma}}_{12}^2 - 2\tilde{\tilde{\sigma}}_{11,12}$  and

$$\tilde{\tilde{\sigma}}_{11,12} := cov\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(1),\phi(1)}), h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(1),\phi(2)})\right) = \sigma_{11,12}\left[\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(1)\phi(1)})\right)'\right]\left[\left(h^{-1}(\widehat{\boldsymbol{\theta}}_{\phi(1),\phi(2)})\right)'\right]A.77\right]$$

A.4. **Proof of THEOREM 2.** In the semi-sparse regime the proof follows the same steps as the proof of Theorem 1. The only difference is in the bias terms, variances, and covariance terms, that are computed according to the following lemma.

**LEMMA A.3.** (Corollary 2 extension for semi-sparse case in Tang et al. (2017)). Let  $\mathbf{A} \sim SBM(\boldsymbol{\xi}, \boldsymbol{\theta}, \rho_n)$  be a  $\tilde{K}$ -block stochastic blockmodel adjacency matrix on n vertices with sparsity factor  $\rho_n$ . Let  $\mu_1, \ldots, \mu_{\tilde{K}}$  be the center of the blocks in  $\mathbb{R}^d$  and let  $\boldsymbol{\theta}_{k\ell} = \boldsymbol{\mu}_k^T \mathbf{I}_{d_1,d_2} \boldsymbol{\mu}_\ell$  be the probability of a link between nodes in blocks k and  $\ell$ . Let K be known and let  $\boldsymbol{\hat{\xi}} : [n] \to [\tilde{K}]$  be the function that assigns nodes to clusters, estimated using GMM or K-means clustering on the rows of  $\boldsymbol{\hat{Y}} = \boldsymbol{\hat{U}} | \boldsymbol{\hat{S}} |^{1/2}$  (as in the proof of Theorem A.1). Let  $\boldsymbol{\hat{\theta}}_{k\ell} = \boldsymbol{\hat{\mu}}_k^T \mathbf{I}_{d_1,d_2} \boldsymbol{\hat{\mu}}_\ell$  and let  $\boldsymbol{\hat{\Delta}} = \sum_{k=1}^K \hat{\eta}_k \boldsymbol{\hat{\mu}}_k \boldsymbol{\hat{\mu}}_\ell^T$ . Define  $\boldsymbol{\Delta} = \sum_{k=1}^K \eta_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T$  and let  $\zeta_{k\ell} = \boldsymbol{\mu}_k^T \boldsymbol{\Delta}^{-1} \boldsymbol{\mu}_\ell$ . Define  $\tilde{\sigma}_{kk}^2$  for  $k \in [K]$  to be

$$\widetilde{\sigma}_{kk}^{2} = 4\theta_{kk}\zeta_{kk}^{2} + 4\sum_{r=1}^{K}\eta_{r}\theta_{kr}\zeta_{kr}^{2}\left(\frac{1}{\eta_{k}} - 2\zeta_{kk}\right)^{2} + 2\sum_{r=1}^{K}\sum_{s=1}^{K}\eta_{r}\eta_{s}\theta_{rs}\zeta_{kr}^{2}\zeta_{ks}^{2} \quad (A.78)$$

and define  $\widetilde{\sigma}_{k\ell}^2$  for  $k \in [\tilde{K}]$  and  $\ell \in [\tilde{K}]$ ,  $k \neq \ell$  to be

$$\widetilde{\sigma}_{k\ell}^2 = (\boldsymbol{\theta}_{kk} + \boldsymbol{\theta}_{\ell\ell}) \zeta_{kk}^2 + 2\boldsymbol{\theta}_{k\ell} \zeta_{kk} \zeta_{\ell\ell} + \sum_{r=1}^K \eta_r \boldsymbol{\theta}_{kr} \zeta_{\ell r}^2 \left(\frac{1}{\eta_k} - 2\zeta_{kk}\right)$$
(A.79)

+ 
$$\sum_{r=1}^{K} \eta_r \boldsymbol{\theta}_{\ell r} \zeta_{kr}^2 \left( \frac{1}{\eta_{\ell}} - 2\zeta_{\ell} \right) - 2 \sum_{r=1}^{K} \eta_r \left( \boldsymbol{\theta}_{kr} + \boldsymbol{\theta}_{\ell r} \right) \zeta_{kr} \zeta_{r\ell} \zeta_{k\ell}$$
(A.80)

+ 
$$\frac{1}{2} \sum_{r=1}^{K} \sum_{s=1}^{K} \eta_r \eta_s \boldsymbol{\theta}_{rs} \left( \zeta_{kr} \zeta_{\ell s} + \zeta_{\ell r} \zeta_{ks} \right)^2$$
. (A.81)

Let  $\ddot{\psi}_{k\ell}$  be defined as

$$\ddot{\psi}_{k\ell} = \sum_{r=1}^{K} \widehat{\eta}_r \left( \widehat{\theta}_{kr} + \widehat{\theta}_{\ell r} \right) \widehat{\mu}_k^T \widehat{\Delta}^{-1} \widehat{\mu}_\ell$$
(A.82)

$$-\sum_{r=1}^{K}\sum_{s=1}^{K}\widehat{\mu}_{k}\widehat{\mu}_{s}\widehat{\theta}_{sr}\widehat{\mu}_{s}^{T}\widehat{\Delta}^{-1}\mathbf{I}_{d_{1},d_{2}}\widehat{\Delta}^{-1}\left(\widehat{\mu}_{\ell}\widehat{\mu}_{k}^{T}+\widehat{\mu}_{l}\widehat{\mu}_{\ell}^{T}\right)\widehat{\Delta}^{-1}\widehat{\mu}_{s}.$$
 (A.83)

Then there exists a (sequence of) permutation(s)  $\phi \equiv \phi_n$  on [K] such that for any  $k \in [K]$ and  $\ell \in [K]$ ,

$$n\rho_n^{1/2} \left( \widehat{\boldsymbol{\theta}}_{\phi(k),\phi(\ell)} - \boldsymbol{\theta}_{k\ell} - \frac{\ddot{\psi}_{k\ell}}{n\rho_n} \right) \stackrel{d}{\to} N(0,\widetilde{\sigma}_{k\ell}) \tag{A.84}$$

as  $n \to \infty$ ,  $\rho_n \to 0$ , and  $n\rho_n = \omega(\sqrt{n})$ .

Carey Business School, Johns Hopkins University, 100 International Dr., Baltimore, MD 21202

Department of Sociology, Johns Hopkins University, 3400 North Charles St., Baltimore, MD 21218

Department of Statistics, University of Michigan, 1085 South University Ave., Ann Arbor, MI 48109

DEPARTMENT OF APPLIED MATHEMATICS AND STATISTICS, JOHNS HOPKINS UNIVERSITY, 3400 N CHARLES ST., BALTIMORE, MD 21218