# Community Detection in High-Dimensional Graph Ensembles

Robert Malinas, Dogyoon Song, Alfred O. Hero III

Department of Electrical & Computer Engineering University of Michigan

Ann Arbor, MI, 48105, USA

Email: {rmalinas, dogyoons, hero}@umich.edu

Abstract—Detecting communities in high-dimensional graphs can be achieved by applying random matrix theory where the adjacency matrix of the graph is modeled by a Stochastic Block Model (SBM). However, the SBM makes an unrealistic assumption that the edge probabilities are homogeneous within communities, i.e., the edges occur with the same probabilities. The Degree-Corrected SBM is a generalization of the SBM that allows these edge probabilities to be different, but existing results from random matrix theory are not directly applicable to this heterogeneous model. In this paper, we derive a transformation of the adjacency matrix that eliminates this heterogeneity and preserves the relevant eigenstructure for community detection. We propose a test based on the extreme eigenvalues of this transformed matrix and (1) provide a method for controlling the significance level, (2) formulate a conjecture that the test achieves power one for all positive significance levels in the limit as the number of nodes approaches infinity, and (3) provide empirical evidence and theory supporting these claims.

Index Terms—community detection, random matrix theory, degree-corrected stochastic block model

#### I. INTRODUCTION

Networks are composed of nodes and pairs of nodes, known as edges or connections, and are capable of representing a diverse range of data. For instance, social networks consist of nodes representing users and edges denoting interpersonal relationships [47]. In satellite communication networks, satellites and microwave channels comprise the nodes and edges, respectively [15]. Due to their versatility, statistical inference with network data is important to a variety of scientific studies such as social science [43], metabolism [19], and epidemiology [30] (see [32] for a review).

Networks often contain communities—groups of nodes within which connections are especially dense [34]. Determining whether or not a network contains communities, which is known as community detection, is a fundamental problem in network data analysis [17]. Methods for community detection typically fall into one of three categories [48]: greedy algorithms [34], global optimization methods [33], [36], [39], or probabilistic models [1], [2], [18], [24], [35].

In this work, we address community detection within a probabilistic model using the framework of statistical hypothesis testing. The Erdős–Rényi Model (ERM) [16] is perhaps

the simplest probabilistic network model, in which n nodes are connected independently at random with some probability  $p \in (0, 1)$ ; however, modeling multiple communities requires additional degrees of freedom. A generalization of the ERM, the Stochastic Block Model (SBM) [18] consists of n nodes, each of which belongs to one of  $1 \le k \le n$  communities, a community assignment function  $\varphi : [n] \to [k]$ , and an irreducible symmetric matrix of probabilities  $P \in (0, 1)^{k \times k}$ . One typically thinks of k as being much smaller than n. A network is drawn from the SBM if any pair of nodes  $i, j \in [n]$  is connected with probability  $(P)_{\varphi(i)\varphi(j)} = (P)_{\varphi(j)\varphi(i)}$  mutually independently from all other node pairs. Although simple, the SBM provides a capable sandbox for describing many interesting problems, e.g., planted partition (or clique) detection [14], and exhibits interesting phase transitions (see [1] for a recent survey). The hypotheses for community detection in SBMs is straightforward: the null hypothesis is that k = 1, i.e., the SBM is an ERM. The alternative hypothesis is that k > 1.

Despite its merits, the SBM has limitations in capturing real-world network features like the rich-club phenomenon, where high-degree nodes tend to connect, even across different communities [13], [49]. More generally, the SBM is too simplistic to represent high levels of degree variation within a single community [27], leading to underfitting in many real networks. This underfitting increases the type I error rate, i.e., the chance of incorrectly declaring that more than one community is present, in the context of statistical hypothesis testing for community detection [22].

To remedy this underfitting, the authors of [24] proposed the Degree-Corrected Stochastic Block Model (DCSBM). This generalization of the SBM includes additional parameters to capture node-specific connection affinity: the likelihood for a node to connect to other nodes regardless of community membership (cf. Definition II.1). However, while offering improved modeling capabilities, working with this more complex model presents significant mathematical challenges.

This paper proposes a statistical test, based on random matrix theory, for community detection within a large DCSBM. Our key contributions are summarized here:

 We introduce a transformation method that converts the heterogeneous and intractable adjacency matrix into a tractable Wigner ensemble. Additionally, we present an

This work was partially supported by Dept of the Air force grant FA8650-19-C-1712 and Army Research Office grant W911NF2310343.

approach for directly estimating the transformed matrix from a single snapshot of the adjacency matrix.

- 2) We propose a method to determine the significance level of the test, which is exact in the limit that n approaches infinity.
- 3) We argue that our test achieves perfect detection for all positive significance levels as n tends to infinity, and provide empirical evidence of this claim.

In Section II, we provide a formal definition of the DCSBM and identify the main problem addressed by this paper. In Section III, we overview related work. In Section IV, we initiate our analysis, providing relevant results from random matrix theory when necessary. In Section V, we propose the test, provide a method for estimating its false alarm rate, and argue for its consistency in the limit as the number of nodes approaches infinity. In Section VI, we provide empirical evidence supporting the assertions of Section V. In Section VII, we conclude and suggest directions for future study.

#### **II. PRELIMINARIES**

We begin with a formal definition of the Degree-Corrected Stochastic Block Model (DCSBM).

Definition II.1 (Degree-Corrected Stochastic Block Model). Let  $n, k \in \mathbb{N}$  such that k < n, and let  $\epsilon \in (0, 1/2]$ . Furthermore, let

- φ: [n] → [k] be surjective;
  W = (w<sub>µν</sub>)<sup>k</sup><sub>µ,ν=1</sub> ∈ ℝ<sup>k×k</sup><sub>+</sub> be symmetric and irreducible;
  θ = (θ<sub>1</sub>,...,θ<sub>n</sub>) ∈ (0,1]<sup>n</sup>;

such that

- $\sum_{i \in \varphi^{-1}(\{\mu\})} \theta_i = 1 \text{ for all } \mu \in [k];$ •
- $w_{\varphi(i)\varphi(j)}\theta_i\theta_j \in [\epsilon, 1-\epsilon]$  for all  $i, j \in [n]$ .

A random matrix  $A \in \{0,1\}^{n \times n}$  is an adjacency matrix drawn from the Degree-Corrected Stochastic Block Model with parameter  $(\epsilon, \varphi, \theta, W)$ , written  $A \sim \text{DCSBM}_n(\epsilon, \varphi, \theta, W)$ , if (*i*) A is symmetric;

- (ii)  $\{A_{ij}\}_{1 \le i \le j \le n}$  is a mutually independent set of random variables;
- (iii) for all  $i, j \in [n]$ , if  $i \in \varphi^{-1}(\{\mu\})$  and  $j \in \varphi^{-1}(\{\nu\})$ , then<sup>1</sup>  $(A)_{ij} \sim \text{Bern} (\theta_i \theta_j w_{\mu\nu}).$

The symbols n and k will be used throughout the rest of this paper to refer to the number of nodes and the number of communities, respectively, in a DCSBM model. Interpretations of the parameters  $\theta$  and W may be found in [24]. We aim to test the hypotheses

$$H_0: k = 1,$$
  
 $H_1: k > 1.$ 

The crux of this paper is that the model under  $H_0$  is not, in general, an ERM; rather, it is a DCSBM with 1 community. The expected degree distribution of a one-community DCSBM may be arbitrary, whereas the ERM expected degree distribution is always a single atom with mass one.

## III. RELATED WORK

We overview related work on community detection methods that are based on matrices such as the adjacency, modularity, or graph Laplacian matrices, within a DCSBM framework.

a) Signed Polygon Statistics: A class of Signed Polygon statistics proposed in [21] assigns scores to each m-gon in a network for some  $m \geq 3$ . The scores are based on the degree of each node in the m-gon, and the statistic is the sum of all such scores. The detection performance of Signed Polgyon statistics is shown to be robust to sparsity and mixed membership, in which communities may overlap. Furthermore, these statistics are capable of detecting small planted cliques with size on the order  $n^{1/2}$  [22].

b) Spectral Clustering: The authors of [20] consider entrywise ratios of a small number of leading eigenvectors of the adjacency matrix. It is argued that these ratio vectors effectively mod out community-independent degree heterogeneity, and clustering the (k-1)-tuple of ratios for each node via k-means is capable of detecting communities. The authors of [5] characterize the empirical spectral distributions of a class of normalized modularity matrices drawn from the DCSBM. In addition to describing phase transitions for community detection, they propose a spectral clustering algorithm that finds an optimal normalization for the modularity matrix. followed by a k-means clustering of its eigenvectors.

c) Extreme Eigenvalue Tests for Special DCSBMs: Perhaps the first work that uses the extreme (the largest and smallest) eigenvalues for community detection, [31] establishes a phase transition in the largest eigenvalue of the modularity matrix for a special case of the DCSBM, namely the planted partition model. The authors of [8] also test the extreme eigenvalues of a transformed adjacency matrix. The transformation is similar to that of this paper in that they are both entrywise centerings and rescalings of the adjacency matrix; however, [8] does not consider a general DCSBM null hypothesis, which is the main focus of this paper. Finally, a goodness-of-fit test is proposed in [29] in which  $k = k_0$ vs  $k \neq k_0$  is tested sequentially for  $k_0 = 1, 2, ...$  until the  $k = k_0$  hypothesis is accepted. This method can be adapted to community detection by terminating the sequence after testing  $k_0 = 1$ ; however, like [8], [29] does not characterize the significance level under a DCSBM null hypothesis.

## IV. COMMUNITY DETECTION FOR THE DCSBM USING **RANDOM MATRIX THEORY**

Let  $n, k, \varphi, \theta = (\theta_1, \ldots, \theta_n)$ , and  $W = (w_{\mu\nu})_{\mu,\nu=1}^k$ be as in Definition II.1, and suppose  $A = (a_{ij})_{i,j=1}^n \sim$  $DCSBM_n(\epsilon, \varphi, \theta, W)$ . We may write

$$A = \mathbb{E}A + (A - \mathbb{E}A).$$

A. Analysis Under One-Community DCSBM Null Hypothesis Under the null hypothesis  $H_0$ ,

$$\mathbb{E}A = w_{11} \boldsymbol{\theta} \boldsymbol{\theta}^T$$

hence, A is the sum of a rank-one matrix  $\mathbb{E}A$  and a centered random matrix  $A - \mathbb{E}A$ . Moreover, the entries in the diagonal

<sup>&</sup>lt;sup>1</sup>A random variable X is Bernoulli with parameter  $p, p \in [0, 1]$ , which we denote by  $X \sim \text{Bern}(p)$ , if X = 1 with probability p and X = 0 otherwise.

and upper triangle  $\{a_{ij} : 1 \leq i \leq j \leq n\}$  are mutually independent. The difficulty with directly analyzing A under the DCSBM is that the variances

$$s_{ij} \coloneqq \mathbb{E} |a_{ij} - \mathbb{E} a_{ij}|^2 = w_{11} \theta_i \theta_j (1 - w_{11} \theta_i \theta_j), \quad \forall i, j \in [n],$$

are, in general, heterogeneous in that they depend on i, j. Due to this heterogeneity, the eigenvalues of  $A - \mathbb{E}A$  can only be described implicitly. Specifically, the empirical spectral distribution of  $A - \mathbb{E}A$  (cf. Definition IV.5) is asymptotically characterized via its Stieltjes transform as the implicit solution to a quadratic vector equation [3], [4], the fundamental properties of which have only recently been described. Even less is known about finite-rank perturbations of such a matrix, making it difficult to describe the extreme eigenvalues of A.

## B. Special Case with Homogeneous Parameters

In the special case of homogeneous parameters where  $\theta_i = \theta_j$  for all  $i, j \in [n]$ , the DCSBM model under the null hypothesis reduces to an ERM with parameter  $w_{11}\theta_1^2$ . In this special case,  $A - \mathbb{E}A$  is distributed as a scaled Wigner ensemble [44]–[46], defined below.

**Definition IV.1** (Wigner ensemble [6]). Let  $H = (h_{ij})_{i,j=1}^n \in \mathbb{C}^{n \times n}$  be a Hermitian random matrix. The random matrix ensemble H is a Wigner ensemble if its entries are centered and normalized, i.e.,  $\mathbb{E}h_{ij} = 0$  and  $\mathbb{E}|h_{ij}|^2 = \frac{1}{n}$  for all  $i, j \in [n]$ , and its entries in the diagonal and upper triangle  $\{h_{ij} : 1 \le i \le j \le n\}$  are mutually independent.

Much is known about Wigner ensembles and, in particular, their extreme eigenvalues and the extreme eigenvalues of low-rank perturbations thereof. Therefore, it is desirable to work with Wigner ensembles when possible. To this end, we define a map that transforms a DCSBM adjacency matrix A to a (scaled) Wigner matrix B.

**Proposition IV.2.** Let  $n, k, \varphi, \theta = (\theta_1, \ldots, \theta_n)$ , and  $W = (w_{\mu\nu})_{\mu,\nu=1}^k$  be as in Definition II.1, and suppose  $A \sim \text{DCSBM}_n(\varphi, \theta, W)$ . Furthermore, let  $B \in \mathbb{R}^{n \times n}$  such that

$$(B)_{ij} = \frac{(A)_{ij} - \theta_i \theta_j w}{\sqrt{\theta_i \theta_j w (1 - \theta_i \theta_j w)}}, \qquad \forall i, j \in [n],$$

where  $w \equiv w_{11}$ . If k = 1, then  $n^{-1/2} \cdot B$  is a Wigner ensemble.

*Proof.* It follows from (i) and (ii) of Definition II.1 that B is symmetric with independent entries in the diagonal and upper triangle. The proof is complete by noticing that for any  $i, j \in [n]$ ,

1)

$$\mathbb{E}\left[n^{-1/2} \cdot (B)_{ij}\right] = n^{-1/2} \cdot \frac{\mathbb{E}(A)_{ij} - \theta_i \theta_j w}{\sqrt{\theta_i \theta_j w (1 - \theta_i \theta_j w)}}$$
$$= n^{-1/2} \cdot \frac{\theta_i \theta_j w - \theta_i \theta_j w}{\sqrt{\theta_i \theta_j w (1 - \theta_i \theta_j w)}}$$
$$= 0,$$

and

$$\mathbb{E} \left| n^{-1/2} \cdot (B)_{ij} \right|^2 = \frac{\mathbb{E} \left| (A)_{ij} - \theta_i \theta_j w \right|^2}{n \cdot \theta_i \theta_j w (1 - \theta_i \theta_j w)}$$
$$= \frac{\theta_i \theta_j w - (\theta_i \theta_j w)^2}{n \cdot \theta_i \theta_j w (1 - \theta_i \theta_j w)}$$
$$= \frac{1}{n}.$$

#### C. Properties of Wigner Ensembles

2)

We now discuss some important results from random matrix theory. Following the seminal work of Baik, Ben Arous, and Péché [7], which studies the extreme eigenvalues of 'spiked' sample covariance matrices [23], significant attention was devoted to developing a parallel theory for low-rank perturbations of Wigner ensembles [10], [11], [25], [26], [28], [37], [38]. The key relevant results of this line of work, under mild technical conditions, include:

- 1) the distributions of the extreme eigenvalues;
- a sharpening of the empirical spectral distribution asymptotics (cf. Definition IV.5) from [44]–[46];
- a phase transition for the extreme eigenvalues under an additive low-rank perturbation.

Next, we provide the specific results required for our community detection analysis.

1) Distribution of the Extreme Eigenvalues: Under mild technical conditions, the marginal distributions of the extreme eigenvalues of an appropriately centered-and-rescaled Wigner ensemble each converge weakly (cf. [9, Section 25]) to the Tracy-Widom distribution, defined next.

**Definition IV.3** (Tracy-Widom ensemble [40], [41]). *The Tracy-Widom distribution*  $TW_1$  *is the probability measure on*  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  *with CDF* 

$$F_1(x) = \exp\left[-\frac{1}{2}\int_x^\infty q(y) + (y-x)q^2(y)\mathrm{d}y\right]$$

where q(y) is the unique solution to the Painlevé II differential equation

$$\frac{\mathrm{d}^2}{\mathrm{d}y^2}q(y) = yq(y) + 2q^3(y),$$

satisfying the boundary condition  $q(y) \asymp \operatorname{Ai}(y)$  as  $y \to \infty$ , where  $\operatorname{Ai}(y)$  denotes the Airy function of the first kind.

The Tracy-Widom density, i.e.,  $\frac{dF_1}{dx}$  is plotted in red in Figure 2.

**Theorem IV.4** ( [28, excerpted from Theorem 1.2]). For each  $n \in \mathbb{N}$ , let  $H_n$  be an  $n \times n$  Wigner ensemble. Suppose that

$$\lim_{s \to +\infty} s^4 \cdot P\left( \left| n^{1/2} (H)_{12} \right| \ge s \right) = 0.$$
 (1)

Then

$$P\left(n^{2/3} \cdot (\lambda_1(H_n) - 2) \le x\right) \to F_1(x), \quad \forall x \in \mathbb{R}.$$

Moreover, the logical inverse holds, i.e., (1) is necessary. A similar result holds for  $-\lambda_n(H_n)$ .

Trivially,  $n^{-1/2} \cdot B$  satisfies (1) because its entries are bounded almost surely.

2) Bulk Characterization of the Spectrum: Other key results offer characterizations of the bulk of the spectrum of a Wigner ensemble, i.e., the eigenvalues outside of an epsilon neighborhood of the extreme eigenvalues. Namely, the empirical spectral distribution, defined below, of a Wigner matrix converges almost surely to the semicircle law.

**Definition IV.5** (Empirical spectral distribution). Let  $H \in \mathbb{C}^{n \times n}$  be a Hermitian matrix. The empirical spectral distribution (ESD) of H is the probability measure

$$\mu_H \coloneqq \frac{1}{n} \sum_{i=1}^n \delta_{\lambda_i(H)},\tag{2}$$

where  $(\lambda_i(H))_{i=1}^n$  is the multiset of eigenvalues of H, and  $\delta_x$  is the Dirac measure with support  $\{x\}$  for any  $x \in \mathbb{R}$ .

Define

$$\mu_{sc}(E) = \int_{E} \rho_{sc}(x) \mathrm{d}x, \qquad \forall E \in \mathcal{B}(\mathbb{R}),$$
$$\rho_{sc}(x) \coloneqq \frac{1}{2\pi} \sqrt{4 - x^2} \mathbb{1}_{[-2, 2]} \text{ for all } x \in \mathbb{R}.$$

**Theorem IV.6** ( [6, Theorem 2.5]). For each  $n \in \mathbb{N}$ , suppose that  $H_n$  is an  $n \times n$  Wigner ensemble. Then, with probability one.

$$\mu_{H_n} \to \mu_{sc}$$

weakly.

where

#### D. Estimating the Transformed Matrix B

In section IV, we introduced a transformed adjacency matrix B and proved that it is a Wigner ensemble under the null hypothesis; however, B depends on the unknown quantities  $\theta$  and  $w_{11}$ . Therefore, we must estimate B to form a viable statistical test. We propose the estimator

$$(\hat{B})_{ij} = \frac{(A)_{ij} - \frac{(A\mathbf{1}\mathbf{1}^T A)_{ij}}{\mathbf{1}^T A\mathbf{1}}}{\sqrt{\frac{(A\mathbf{1}\mathbf{1}^T A)_{ij}}{\mathbf{1}^T A\mathbf{1}}} \left(1 - \frac{(A\mathbf{1}\mathbf{1}^T A)_{ij}}{\mathbf{1}^T A\mathbf{1}}\right)}, \qquad \forall i, j \in [n]$$

informed by the following theorem.

**Theorem IV.7.** Let  $n, k, \varphi, \theta = (\theta_1, \ldots, \theta_n)$ , and  $W = (w_{\mu\nu})_{\mu,\nu=1}^k$  be as in Definition II.1, and suppose  $A = (a_{ij})_{i,j=1}^n \sim \text{DCSBM}_n(\epsilon, \varphi, \theta, W)$ . If k = 1, then

$$P\left(\left|\left(\frac{A\mathbf{1}\mathbf{1}^{T}A}{\mathbf{1}^{T}A\mathbf{1}}\right)_{ij} - w\theta_{i}\theta_{j}\right| \leq \frac{8}{\epsilon} \left(\sqrt{t}n^{-1/2} + tn^{-1}\right)\right)$$
$$\geq 1 - 2\left(2e^{-2t^{2}} + e^{-\frac{n^{2}\epsilon^{2}}{18}}\right), \qquad \forall t > 0, \ \forall i, j \in [n],$$

where  $w \equiv w_{11}$ .

*Proof.* See Appendix (Section VIII-B).

To summarize,  $\frac{(A\mathbf{1}\mathbf{1}^T A)_{ij}}{\mathbf{1}^T A\mathbf{1}}$  fluctuates around the unknown parameters  $w_{11}\theta_i\theta_j$  at the scale  $n^{-1/2}$  for any  $i, j \in [n]$ . Moreover, the fluctuation is subgaussian. From these facts, we conjecture the following.

**Conjecture 1.** Let  $n, k, \varphi, \theta = (\theta_1, \ldots, \theta_n)$ , and  $W = (w_{\mu\nu})_{\mu,\nu=1}^k$  be as in Definition II.1, and suppose  $A = (a_{ij})_{i,j=1}^n \sim \text{DCSBM}_n(\epsilon, \varphi, \theta, W)$ . Then if k = 1,

$$\left| P\left( n^{2/3} \cdot \left( \lambda_1 \left( n^{-1/2} \cdot \hat{B} \right) - 2 \right) \le x \right) - F_1(x) \right| = o(1),$$

and

$$\left| P\left( n^{2/3} \cdot \left( -\lambda_n \left( n^{-1/2} \cdot \hat{B} \right) - 2 \right) \le x \right) - F_1(x) \right| = o(1),$$
  
for all  $x \in \mathbb{R}$ .

In other words, the distributions of the extreme eigenvalues of  $\hat{B}$  are the same as those of a Wigner ensemble asymptotically.

#### V. TEST STATISTIC

We propose the statistic

$$T = \max\left\{n^{2/3} \cdot \left(\lambda_1 \left(n^{-1/2} \cdot \hat{B}\right) - 2\right)\right\}$$
$$n^{2/3} \cdot \left(-\lambda_n \left(n^{-1/2} \cdot \hat{B}\right) - 2\right)\right\}$$

This statistic is reasonable because

- 1) under the null hypothesis, given Conjecture 1, T fluctuates around 0, following a TW<sub>1</sub> distribution asymptotically;
- 2) under the alternative hypothesis, we expect that  $\begin{vmatrix} \lambda_1 \left( n^{-1/2} \cdot \hat{B} \right) 2 \end{vmatrix} \gg n^{-2/3}$  for a large class of alternative models  $W, \theta$ , and  $\varphi$ , with a similar statement holding for  $-\lambda_n \left( n^{-1/2} \cdot \hat{B} \right)$ .

The latter point comes from a BBP-type phase transition for low-rank perturbations of Wigner ensembles, e.g., [25, Theorem 2.7]. In section VI, we provide empirical evidence that indicates a threshold test based on T attains power one in the limit as  $n \to \infty$  for all positive significance levels, i.e., the test is asymptotically consistent. We leave a thorough power analysis for future work.

#### A. False Alarm Rate

Based on Conjecture 1, we pose the following conjecture on the false alarm rate of the proposed statistic T.

**Conjecture 2.** Let  $G_1 = 1 - F_1$ . If k = 1, then

$$P\left(T \ge G_1^{-1}(\alpha/2)\right) \lesssim \alpha, \quad \forall \alpha \in (0,1).$$

*Remark* 1. We note that  $G_1^{-1}$  exists because  $F_1$  is strictly monotonic: the TW<sub>1</sub> distribution is absolutely continuous.

In summary, rejecting the null hypothesis if  $T \ge G_1^{-1}(\alpha/2)$ yields a type I error of at most  $\alpha$ , for any  $\alpha \in (0,1)$ , at least asymptotically. The argument  $\alpha/2$  of the quantile function  $G_1^{-1}$  is due to Bonferonni correction, because we are simultaneously testing two eigenvalues.

## VI. EMPIRICAL RESULTS

In this section, we present empirical evidence supporting the assertion that the eigenvalues of  $n^{-1/2} \cdot \hat{B}$  behave like those of a Wigner ensemble. For all simulations, we generate  $\theta$  via

$$\theta_i = \frac{X_i}{\sum_{j=1}^n X_j}, \qquad \forall i \in [n],$$

where  $X_j \sim \text{Unif}[0.1, 0.9]$ , for all  $j \in [n]$ . For Figure 3, we set

$$W = D \cdot \begin{pmatrix} 0.4 & 0.2 & 0.2 \\ 0.2 & 0.6 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix} \cdot D,$$

where

$$D = \text{diag}\left(\sum_{j \in \varphi - 1(\{1\})}^{n} X_j, \sum_{j \in \varphi - 1(\{2\})}^{n} X_j, \sum_{j \in \varphi - 1(\{3\})}^{n} X_j\right),$$
  
and

$$\left|\varphi^{-1}(\{\mu\})\right| = \frac{n}{3}, \qquad \forall \mu \in [3]$$

For Figures 1 and 2, we set  $w_{11} = \frac{1}{2} \sum_{j=1}^{n} X_j$ . Figures 1 and 2 support that  $n^{-1/2} \cdot \hat{B}$  has a spectrum close

Figures 1 and 2 support that  $n^{-1/2} \cdot B$  has a spectrum close to that of a Wigner ensemble. In particular, Figure 2 indicates that the extreme eigenvalues of  $n^{-1/2} \cdot \hat{B}$  converge to those of a Wigner ensemble at a rate faster than  $n^{-2/3}$ . Figure 3 offers a glimpse into the asymptotic power of the proposed test, indicating that it is asymptotically one for any positive significance level.



Fig. 1. A histogram of the eigenvalues of a single realization of  $n^{-1/2} \cdot \hat{B}$ under the null hypothesis k = 1, with n = 3000. Overlaid on the histogram is the semicircle density  $\rho_{sc}(x) = \frac{1}{2\pi}\sqrt{4-x^2} \cdot \mathbb{I}_{[-2, 2]}$ .

#### VII. CONCLUSION

In this paper, we proposed a test for detecting communities within a Degree-Corrected Stochastic Block Model. The test is based on the extreme eigenvalues of a an element-wise centered and rescaled adjacency matrix. Roughly, the proposed centering and rescaling are consistent with a transformation that maps the adjacency matrix A to a Wigner ensemble.



Fig. 2. A histogram of  $n^{2/3} \cdot \left(\lambda_1 \left(n^{-1/2} \cdot \hat{B}\right) - 2\right)$  over 2000 independent realizations with n = 500. Overlaid on the histogram is the TW<sub>1</sub> density  $\frac{\mathrm{d}F_1}{\mathrm{d}x}$ , computed in software using [12].



Fig. 3. Receiver operating characteristics for different values of n.

Because of this, we are able to approximate the distribution of the proposed statistic using the Tracy-Widom distribution: the asymptotic distribution of the extreme eigenvalues of a Wigner ensemble. Additionally, we provided a method for controlling the false alarm rate of the proposed test. Future work includes an analysis of the power of this test, which we believe converges as  $n \to \infty$  to one for any positive significance level, and rigorous proofs of Conjectures 1 and 2.

#### VIII. APPENDIX

## A. Notation and Conventions

Let  $\mathbb{N}$  be the set of positive integers. For  $n \in \mathbb{N}$ , let  $[n] := \{1, \ldots, n\}$ . Let  $\mathbb{C}$  be the set of complex numbers, let  $\mathbb{R}$  be the set of real numbers, and define  $\mathbb{R}_+ := (0, \infty]$ . For a matrix  $A \in \mathbb{C}^{m \times n}$ , we write  $(A)_{ij}$  for the element in the  $i^{th}$  row and  $j^{th}$  column of A. Let  $\mathbf{1}_n := (1, 1, \ldots, 1) \in \mathbb{R}^n$ . For  $i, j \in [n]$ , we set  $\delta_{ij} = 1$  if i = j and  $\delta_{ij} = 0$  otherwise. For  $n \in \mathbb{N}$ , we reserve the letter  $I_n \in \mathbb{C}^{n \times n}$  to denote the identity matrix throughout this paper. For a diagonalizable matrix  $A \in \mathbb{C}^n$ 

 $\mathbb{C}^{n \times n}$ , we write  $(\lambda_i(A))_{i=1}^n$  for the multiset of eigenvalues of A such that  $\lambda_1(A) \geq \lambda_2(A) \geq \cdots \geq \lambda_n(A)$ . For a set  $S \subseteq \mathbb{C}$ , we use the notation  $A = (a_{ij})_{i,j=1}^n \in S^{n \times n}$  to denote the matrix A with elements  $(A)_{ij} = a_{ij} \in S$  for all  $i, j \in$ [n]. For a vector  $\mathbf{u} = (u_1, \ldots, u_n) \in \mathbb{C}^n$ , we let diag $(\mathbf{u}) \in$  $\mathbb{C}^{n \times n}$  denote the matrix with elements  $(\text{diag}(\mathbf{u}))_{ij} = \delta_{ij}u_i$ . For a topological space  $\mathcal{X}$ , we use  $\mathcal{B}(\mathcal{X})$  to denote the Borel  $\sigma$ -algebra generated by the open sets in  $\mathcal{X}$ .

## B. Proof of Theorem IV.7

1) Useful Lemmas:

**Theorem VIII.1** (Hoeffding's inequality for bounded random variables [42, Theorem 2.2.6]). Let  $X_1, \ldots, X_n$  be independent random variables. Assume that  $X_i \in [m_i, M_i]$  for every  $i \in [n]$ . Then, for any t > 0, we have

$$P\left(\sum_{i=1}^{n} \left(X_i - \mathbb{E}X_i\right) \ge t\right) \le \exp\left(-\frac{2t^2}{\sum_{i=1}^{n} \left(M_i - m_i\right)^2}\right).$$

**Lemma VIII.2.** Let  $a, b, c, d \in \mathbb{R}$  such that  $b \neq 0$ ,  $b + d \neq 0$ , and  $\left|\frac{a}{b}\right| \leq 1$ . Then

$$\left|\frac{a+c}{b+d} - \frac{a}{b}\right| \le \frac{|c|+|d|}{|b+d|}.$$

Proof.

$$\left|\frac{a+c}{b+d} - \frac{a}{b}\right| = \left|\frac{bc-ad}{b(b+d)}\right|$$
$$\leq \left|\frac{c}{b+d}\right| + \left|\frac{a}{b}\right| \left|\frac{d}{b+d}\right|$$
$$\leq \frac{|c|+|d|}{|b+d|}.$$

### 2) Proof of Main Result:

*Proof.* We seek to apply Lemma VIII.2 with  $a = w^2 \theta_i \theta_j$ , b = w,  $c = (A\mathbf{1}\mathbf{1}^T A)_{ij} - w^2 \theta_i \theta_j$ , and  $d = \mathbf{1}^T A\mathbf{1} - w$ . We begin by establishing fundamental bounds on the absolute "errors"  $|(A\mathbf{1}\mathbf{1}^T A)_{ij} - w^2 \theta_i \theta_j|$  and  $|\mathbf{1}^T A\mathbf{1} - w|$ . Then, we show that these error bounds are "small" via a concentration inequality for sums of independent random variables. Finally, we apply Lemma VIII.2 to yield the result.

Fix  $i, j \in [n]$  and write  $E_{ij} \equiv \left(\frac{A\mathbf{1}\mathbf{1}^T A}{\mathbf{1}^T A\mathbf{1}}\right)_{ij} - w\theta_i\theta_j$ . To this end, we begin with fundamental bounds on the "error" terms  $\mathbf{1}^T A\mathbf{1} - w$  and  $(A\mathbf{1}\mathbf{1}^T A)_{ij} - w^2\theta_i\theta_j$ . We have

$$(A\mathbf{1}\mathbf{1}^T A)_{ij} = \sum_{l=1}^n a_{il} \sum_{m=1}^n a_{mj},$$

and

$$\mathbf{1}^{T} A \mathbf{1} = \sum_{p=1}^{n} \sum_{q=1}^{n} a_{pq}$$
$$= 2 \sum_{p=1}^{n} \sum_{q=p+1}^{n} a_{pq} + \sum_{r=1}^{n} a_{rr}.$$

Then,

$$\left| (A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j} \right| = \left| \sum_{l=1}^{n} a_{il} \sum_{m=1}^{n} a_{mj} - w^{2}\theta_{i}\theta_{j} \right|$$
$$= \left| \left( \sum_{l=1}^{n} a_{il} - w\theta_{i} \right) \left( \sum_{m=1}^{n} a_{mj} - w\theta_{j} \right) \right|$$
$$-2w^{2}\theta_{i}\theta_{j} + w\theta_{j} \sum_{l=1}^{n} a_{il} + w\theta_{i} \sum_{m=1}^{n} a_{mj} \right|$$
$$\leq \left| \sum_{l=1}^{n} a_{il} - w\theta_{i} \right| \left| \sum_{m=1}^{n} a_{mj} - w\theta_{j} \right|$$
$$+ w\theta_{j} \left| \sum_{l=1}^{n} a_{il} - w\theta_{i} \right| + w\theta_{i} \left| \sum_{m=1}^{n} a_{mj} - w\theta_{j} \right|, \qquad (3)$$

and

(

$$\left|\mathbf{1}^{T}A\mathbf{1} - \sum_{p=1}^{n}\sum_{q=1}^{n}w\theta_{p}\theta_{q}\right| \leq \left|\sum_{r=1}^{n}a_{rr} - \sum_{r=1}^{n}\theta_{r}^{2}w\right|$$
(4)

$$+2\left|\sum_{p=1}^{n}\sum_{q=p+1}^{n}a_{pq}-\sum_{p=1}^{n}\sum_{q=p+1}^{n}\theta_{p}\theta_{q}w\right|.$$
(5)

By Theorem VIII.1 (Hoeffding's inequality), for any  $t \ge 0$  we have

(i) 
$$P\left(\left|\sum_{l=1}^{n} a_{il} - \theta_{i}w\right| \ge t\right) \le \exp\left(-\frac{2t^{2}}{n}\right);$$
  
(ii)  $P\left(\left|\sum_{m=1}^{n} a_{mj} - \theta_{j}w\right| \ge t\right) \le \exp\left(-\frac{2t^{2}}{n}\right);$   
(iii)

$$P\left(\left|\sum_{p=1}^{n}\sum_{q=p+1}^{n}a_{pq}-\sum_{p=1}^{n}\sum_{q=p+1}^{n}\theta_{p}\theta_{q}w\right| \ge t\right)$$
$$\le \exp\left(-\frac{2t^{2}}{n^{2}-n}\right) \le \exp\left(-\frac{2t^{2}}{n^{2}}\right);$$

(iv) 
$$P\left(\left|\sum_{r=1}^{n} a_{rr} - \sum_{r=1}^{n} \theta_r^2 w\right| \ge t\right) \le \exp\left(-\frac{2t^2}{n}\right)$$
,  
where in (i) and (ii) we used the normalization co

where in (i) and (ii) we used the normalization condition  $\sum_{l=1}^{n} \theta_l = 1$  (cf. Definition II.1). Combining (3), (i), and (ii), we find from union bound and DeMorgan's law that

$$P\left(\left|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}\right| \leq t^{2} + w(\theta_{i} + \theta_{j})t\right)$$
  
$$\geq 1 - 2\exp\left(-\frac{2t^{2}}{n}\right), \quad \forall t \geq 0.$$
(6)

Similarly, from (iii), (iv), and (4), it follows that

$$P\left(\left|\mathbf{1}^{T} A \mathbf{1} - w\right| \le 3t\right) \ge 1 - \left(\exp\left(-\frac{2t^{2}}{n^{2}}\right) + \exp\left(-\frac{2t^{2}}{n}\right)\right)$$
$$\ge 1 - 2\exp\left(-\frac{2t^{2}}{n^{2}}\right), \quad \forall t \ge 0.$$
(7)

Note that  $w = \sum_{p=1}^{n} \sum_{q=1}^{n} w \theta_p \theta_q \ge n^2 \epsilon$ , thus, (7) implies

$$P\left(\left|\mathbf{1}^{T}A\mathbf{1}-w\right| \leq \frac{w}{2}\right) \geq 1 - 2\exp\left(-\frac{w^{2}}{18n^{2}}\right)$$
$$\geq 1 - 2\exp\left(-\frac{n^{2}\epsilon^{2}}{18}\right). \quad (8)$$

Assuming  $|w - \mathbf{1}^T A \mathbf{1}| \leq \frac{w}{2}$  and, thus,  $\mathbf{1}^T A \mathbf{1} > 0$ , Lemma VIII.2<sup>2</sup> yields

$$E_{ij}| \leq \frac{|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}| + |\mathbf{1}^{T}A\mathbf{1} - w|}{\mathbf{1}^{T}A\mathbf{1}} \\ = \frac{|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}| + |\mathbf{1}^{T}A\mathbf{1} - w|}{|w - (w - \mathbf{1}^{T}A\mathbf{1})|} \\ \leq \frac{|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}| + |\mathbf{1}^{T}A\mathbf{1} - w|}{|w - |w - \mathbf{1}^{T}A\mathbf{1}||} \\ \leq 2\frac{|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}| + |\mathbf{1}^{T}A\mathbf{1} - w|}{|w|} \\ \leq 2\frac{|(A\mathbf{1}\mathbf{1}^{T}A)_{ij} - w^{2}\theta_{i}\theta_{j}| + |\mathbf{1}^{T}A\mathbf{1} - w|}{n^{2}\epsilon}.$$
(9)

Noting that  $w\theta_p = w\theta_p \sum_{q=1}^n \theta_q \le n(1-\epsilon)$  for all  $p \in [n]$  and letting t > 0 be arbitrary, it follows that

$$P\left(|E_{ij}| \le 2\frac{\sqrt{t}(2n^{3/2}(1-\epsilon)) + 4tn}{n^2\epsilon}\right)$$
$$\ge P\left(|E_{ij}| \le 2\frac{\left(\sqrt{tn}\right)^2 + w(\theta_i + \theta_j)\sqrt{tn} + 3tn}{n^2\epsilon}\right)$$
$$\ge 1 - 2\left(2e^{-2t^2} + e^{-\frac{n^2\epsilon^2}{18}}\right),$$

where in the last step we combined (6), (7), (8), and (9). Finally, for any t > 0,

$$2\frac{\sqrt{t}(2n^{3/2}(1-\epsilon)) + 4tn}{n^{2}\epsilon} = \frac{4(1-\epsilon)\sqrt{t}n^{-1/2} + 8tn^{-1}}{\epsilon} \\ \leq \frac{8}{\epsilon} \left(\sqrt{t}n^{-1/2} + tn^{-1}\right).$$

#### REFERENCES

- Emmanuel Abbe. Community detection and stochastic block models: Recent developments. Journal of Machine Learning Research, 18(177):1–86, 2018.
- [2] Edo M Airoldi, David Blei, Stephen Fienberg, and Eric Xing. Mixed membership stochastic blockmodels. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, <u>Advances in Neural Information</u> Processing Systems, volume 21. Curran Associates, Inc., 2008.
- [3] Oskari Ajanki, Laszlo Erdos, and Torben Krüger. Universality for general wigner-type matrices. <u>Probab. Theory Relat. Fields</u>, 169:667–727, 2017.
- [4] Oskari Ajanki, László Erdős, and Torben Krüger. Quadratic vector equations on complex upper half-plane. <u>Memoirs of the American</u> <u>Mathematical Society</u>, 261(1261):0–0, sep 2019.
- [5] Hafiz Tiomoko Ali and Romain Couillet. Improved spectral community detection in large heterogeneous networks. <u>Journal of Machine Learning</u> <u>Research</u>, 18(225):1–49, 2018.
- [6] Z. Bai and Jack Silverstein. <u>Spectral Analysis of Large Dimensional</u> <u>Random Matrices.</u> 01 2010.

$$^2 \text{with } a = w^2 \theta_i \theta_j, b = w, c = (A \mathbf{1} \mathbf{1}^T A)_{ij} - w^2 \theta_i \theta_j, \text{ and } d = \mathbf{1}^T A \mathbf{1} - w$$

- [7] Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. The Annals of Probability, 33(5):1643 – 1697, 2005.
- [8] Peter J. Bickel and Purnamrita Sarkar. Hypothesis testing for automated community detection in networks. <u>Journal of the Royal Statistical</u> <u>Society. Series B (Statistical Methodology)</u>, 78(1):253–273, 2016.
- [9] Patrick Billingsley. Probability and measure. Wiley, 3rd edition, 1995.
- [10] Alex Bloemendal and Bálint Virág. Limits of spiked random matrices II. <u>The Annals of Probability</u>, 44(4):2726 – 2769, 2016.
- [11] Alex Bloemendal and Bálint Virág. Limits of spiked random matrices i. <u>Probability Theory and Related Fields</u>, 156(3–4):795–825, September 2012.
- [12] Marco Chiani. Distribution of the largest eigenvalue for real wishart and gaussian random matrices and a simple approximation for the tracy-widom distribution. <u>Journal of Multivariate Analysis</u>, 129:69–81, August 2014.
- [13] V. Colizza, A. Flammini, M. A. Serrano, and A. Vespignani. Detecting rich-club ordering in complex networks. <u>Nature Physics</u>, 2(2):110–115, January 2006.
- [14] Anne Condon and Richard M. Karp. Algorithms for graph partitioning on the planted partition model. <u>Random Structures & Algorithms</u>, 18(2):116–140, 2001.
- [15] Bruce R. Elbert. <u>Introduction to satellite communication</u>. Artech House, 2008.
- [16] P. Erdös and A. Rényi. On random graphs i. <u>Publicationes Mathematicae</u> Debrecen, 6:290, 1959.
- [17] Santo Fortunato. Community detection in graphs. <u>Physics Reports</u>, 486(3–5):75–174, February 2010.
- [18] Paul W. Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. <u>Social Networks</u>, 5(2):109–137, 1983.
- [19] H. Jeong, B. Tombor, R. Albert, Z. N. Oltvai, and A.-L. Barabási. The large-scale organization of metabolic networks. <u>Nature</u>, 407(6804):651– 654, Oct 2000.
- [20] Jiashun Jin. Fast community detection by score. <u>The Annals of Statistics</u>, 43(1), February 2015.
- [21] Jiashun Jin, Zheng Tracy Ke, and Shengming Luo. Optimal adaptivity of signed-polygon statistics for network testing, 2019.
- [22] Jiashun Jin, Zheng Tracy Ke, Paxton Turner, and Anru R. Zhang. Phase transition for detecting a small community in a large network, 2023.
- [23] Iain M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. <u>The Annals of Statistics</u>, 29(2):295 – 327, 2001.
- [24] Brian Karrer and M. E. J. Newman. Stochastic blockmodels and community structure in networks. <u>Phys. Rev. E</u>, 83:016107, Jan 2011.
- [25] Antti Knowles and Jun Yin. The isotropic semicircle law and deformation of wigner matrices. <u>Communications on Pure and Applied</u> <u>Mathematics</u>, 66(11):1663–1749, 2013.
- [26] Antti Knowles and Jun Yin. The outliers of a deformed wigner matrix. The Annals of Probability, 42(5):1980–2031, 2014.
- [27] Andrea Lancichinetti, Santo Fortunato, and Filippo Radicchi. Benchmark graphs for testing community detection algorithms. <u>Physical</u> Review E, 78(4), October 2008.
- [28] Ji Oon Lee and Jun Yin. A necessary and sufficient condition for edge universality of wigner matrices. <u>Duke Mathematical Journal</u>, 163(1), January 2014.
- [29] Jing Lei. A goodness-of-fit test for stochastic block models. <u>The Annals</u> of Statistics, 44(1):401 – 424, 2016.
- [30] Cristopher Moore and M. E. J. Newman. Epidemics and percolation in small-world networks. Physical Review E, 61(5):5678–5682, May 2000.
- [31] Raj Rao Nadakuditi and M. E. J. Newman. Graph spectra and the detectability of community structure in networks. <u>Physical Review</u> <u>Letters</u>, 108(18), May 2012.
- [32] M. E. J. Newman. The structure and function of complex networks. SIAM Review, 45(2):167–256, January 2003.
- [33] M. E. J. Newman. Modularity and community structure in networks. <u>Proceedings of the National Academy of Sciences</u>, 103(23):8577–8582, June 2006.
- [34] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. Physical Review E, 69(2), February 2004.
- [35] M. E. J. Newman and E. A. Leicht. Mixture models and exploratory analysis in networks. <u>Proceedings of the National Academy of Sciences</u>, 104(23):9564–9569, June 2007.

- [36] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, <u>Advances in Neural Information Processing Systems</u>, volume 14. MIT Press, 2001.
- [37] S. Péché. The largest eigenvalue of small rank perturbations of hermitian random matrices. <u>Probability Theory and Related Fields</u>, 134(1):127– 173, Jan 2006.
- [38] Alessandro Pizzo, David Renfrew, and Alexander Soshnikov. On finite rank deformations of wigner matrices. <u>Annales de l'Institut Henri</u> Poincaré, Probabilités et Statistiques, 49(1):64–94, Feb 2013.
- [39] Jianbo Shi and J. Malik. Normalized cuts and image segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 22(8):888– 905, 2000.
- [40] Craig A. Tracy and Harold Widom. Level-spacing distributions and the airy kernel. <u>Communications in Mathematical Physics</u>, 159(1):151–174, 1994.
- [41] Craig A. Tracy and Harold Widom. On orthogonal and symplectic matrix ensembles. <u>Communications in Mathematical Physics</u>, 177(3):727–754, 1996.
- [42] Roman Vershynin. <u>High-Dimensional Probability: An Introduction with</u> <u>Applications in Data Science</u>. Number 47 in Cambridge Series in <u>Statistical and Probabilistic Mathematics</u>. Cambridge University Press.
- [43] Stanley Wasserman and Katherine Faust. <u>Social Network Analysis:</u> <u>Methods and Applications</u>. Structural Analysis in the Social Sciences. Cambridge University Press, 1994.
- [44] Eugene P. Wigner. On the statistical distribution of the widths and spacings of nuclear resonance levels. <u>Mathematical Proceedings of the</u> Cambridge Philosophical Society, 47(4):790–798, 1951.
- [45] Eugene P. Wigner. Characteristic vectors of bordered matrices with infinite dimensions. Annals of Mathematics, 62(3):548–564, 1955.
- [46] Eugene P. Wigner. On the distribution of the roots of certain symmetric matrices. The Annals of Mathematics, 67(2):325, 1958.
- [47] Wayne W. Zachary. An information flow model for conflict and fission in small groups. <u>Journal of Anthropological Research</u>, 33(4):452–473, 1977.
- [48] Yunpeng Zhao, Elizaveta Levina, and Ji Zhu. Consistency of community detection in networks under degree-corrected stochastic block models. The Annals of Statistics, 40(4), August 2012.
- [49] Shi Zhou and R.J. Mondragon. The rich-club phenomenon in the internet topology. <u>IEEE Communications Letters</u>, 8(3):180–182, 2004.