

# Edge Label Inference in Generalized Stochastic Block Models: from Spectral Theory to Impossibility Results

Jiaming Xu, Laurent Massoulié, Marc Lelarge

### ▶ To cite this version:

Jiaming Xu, Laurent Massoulié, Marc Lelarge. Edge Label Inference in Generalized Stochastic Block Models: from Spectral Theory to Impossibility Results. Conference on Learning Theory, Jun 2014, Barcelona, Spain. pp.903-920. hal-01066047

## HAL Id: hal-01066047 https://hal.archives-ouvertes.fr/hal-01066047

Submitted on 12 Nov 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Edge Label Inference in Generalized Stochastic Block Models: from Spectral Theory to Impossibility Results

Jiaming Xu

University of Illinois at Urbana-Champaign

Laurent Massoulié Microsoft Research-Inria Joint Centre

Marc Lelarge INRIA-ENS JXU18@ILLINOIS.EDU

LAURENT.MASSOULIE@INRIA.FR

MARC.LELARGE@ENS.FR

#### Abstract

The classical setting of community detection consists of networks exhibiting a clustered structure. To more accurately model real systems we consider a class of networks (i) whose edges may carry labels and (ii) which may lack a clustered structure. Specifically we assume that nodes possess latent attributes drawn from a general compact space and edges between two nodes are randomly generated and labeled according to some unknown distribution as a function of their latent attributes. Our goal is then to infer the edge label distributions from a partially observed network. We propose a computationally efficient spectral algorithm and show it allows for asymptotically correct inference when the average node degree could be as low as logarithmic in the total number of nodes. Conversely, if the average node degree is below a specific constant threshold, we show that no algorithm can achieve better inference than guessing without using the observations. As a byproduct of our analysis, we show that our model provides a general procedure to construct random graph models with a spectrum asymptotic to a pre-specified eigenvalue distribution such as a power-law distribution.

Keywords: Community Detection, Stochastic Blockmodel, Spectral Methods, Galton-Watson Tree

#### 1. Introduction

Detecting communities in networks has received a large amount of attention and has found numerous applications across various disciplines including physics, sociology, biology, statistics, computer science, etc (see the exposition Fortunato (2010) and the references therein). Most previous work assumes networks can be divided into groups of nodes with dense connections internally and sparser connections between groups, and considers random graph models with some underlying cluster structure such as the stochastic blockmodel (SBM), a.k.a. the *planted partition model*. In its simplest form, nodes are partitioned into clusters, and any two nodes are connected by an edge independently at random with probability p if they are in the same cluster and with probability qotherwise. The problem of cluster recovery under the SBM has been extensively studied and many efficient algorithms with provable performance guarantees have been developed (see e.g., Chen and Xu (2014) and the references therein).

Real networks, however, may not display a clustered structure; the goal of community detection should then be redefined. As observed in Heimlicher et al. (2012), interactions in many real networks can be of various types and prediction of unknown interaction types may have practical merit such as prediction of missing ratings in recommender systems. Therefore an intriguing question arises: Can we accurately predict the unknown interaction types in the absence of a clustered structure? To answer it, we generalize the SBM by relaxing the cluster assumption and allowing edges to carry labels. In particular, each node has a latent attribute coming from a general compact space and for any two nodes, an edge is first drawn and then labeled according to some unknown distribution as a function of their latent attributes. Given a partial observation of the labeled graph generated as above, we aim to infer the edge label distributions, which is relevant in many scenarios such as:

- Collaborative filtering: A recommender system can be represented as a labeled bipartite graph where if a user rates a movie, then there is a labeled edge between them with the label being the rating. One would like to predict the missing ratings based on the observation of a few ratings.
- Link type prediction: A social network can be viewed as a labeled graph where if a person knows another person, then there is a labeled edge between them with the label being their relationship type (either friend or colleague). One would like to predict the unknown link types based on the few known link types.
- Prediction of gene expression levels: A DNA microarray can be looked as a a labeled bipartite graph where if a gene is expressed in a sample, then there is a labeled edge between them with the label being the expression level. One would like to predict the unobserved expression level based on the few observed expression levels.

#### 1.1. Problem formulation

The generalized stochastic blockmodel (GSBM) is formally defined by seven parameters n,  $\mathcal{X}$ , P, B,  $\mathcal{L}$ ,  $\mu$ ,  $\omega$ , where n is a positive integer;  $\mathcal{X}$  is a compact space endowed with the probability measure P;  $B : \mathcal{X} \times \mathcal{X} \to [0, 1]$  is a function symmetric in its two arguments;  $\mathcal{L}$  is a finite set with  $\mathcal{P}(\mathcal{L})$  denoting the set of probability measures on it;  $\mu : \mathcal{X} \times \mathcal{X} \to \mathcal{P}(\mathcal{L})$  is a measure-valued function symmetric in its two arguments;  $\omega$  is a positive real number.

**Definition 1** Suppose that there are *n* nodes indexed by  $i \in \{1, ..., n\}$ . Each node *i* has an attribute  $\sigma_i$  drawn in an i.i.d. manner from the distribution *P* on  $\mathcal{X}$ . A random labeled graph is generated based on  $\sigma$ : For each pair of nodes *i*, *j*, independently of all others, we draw an edge between them with probability  $B_{\sigma_i,\sigma_j}$ ; then for each edge (i, j), independently of all others, we label it by  $\ell \in \mathcal{L}$  with probability  $\mu_{\sigma_i,\sigma_j}(\ell)$ ; finally each labeled edge is retained with probability  $\omega/n$  and erased otherwise.

Given a random labeled graph G generated as above, our goal is to infer the edge label distribution  $\mu_{\sigma_i,\sigma_j}$  for any pair of nodes *i* and *j*. To ensure the inference is feasible, we shall make the following *identifiability* assumption: Let  $\nu_{x,y} := B_{x,y}\mu_{x,y}$  and

$$\forall x \neq x' \in \mathcal{X}, \ \sum_{\ell \in \mathcal{L}} \int_{\mathcal{X}} |\nu_{x,y}(\ell) - \nu_{x',y}(\ell)| P(dy) > 0;$$
(1)

otherwise x, x' are statistically indistinguishable and can be combined as a single element in  $\mathcal{X}$ . We emphasize that the model parameters  $(\mathcal{X}, P, B, \mathcal{L}, \mu)$  are all fixed and do not scale with n, while  $\omega$ 

could scale with n. Notice that  $n\omega$  characterizes the total number of observed edge labels and thus can be seen as a measure of "signal strength".

#### 1.2. Main results

We show that it is possible to make meaningful inference of edge label distributions without knowledge of any model parameters in the relatively "sparse" graph regime with  $\omega = \Omega(\log n)$ . In particular, we propose a computationally efficient spectral algorithm with a random weighing strategy. The random weighing strategy assigns a random weight to each label and constructs a weighted adjacency matrix of the label graph. The spectral algorithm embeds the nodes into a finite, low dimensional Euclidean space based on the leading eigenvectors of the weighted adjacency matrix and uses the empirical frequency of labels on the local neighborhood in the Euclidean space to estimate the underlying true label distribution.

In the very "sparse" graph regime with  $\omega = O(1)$ , since there exist at least  $\Theta(n)$  isolated nodes without neighbors and to infer the edge label distribution between two isolated nodes the observed labeled graph G does not provide any useful information, it is impossible to make meaningful inference for at least a positive fraction of node pairs. Moreover, we show that it is impossible to make meaningful inference for any randomly chosen pair of nodes when  $\omega$  is below a specific non-trivial threshold.

As a byproduct of our analysis, we show how the GSBM can generate random graph models with a spectrum asymptotic to a pre-specified eigenvalue distribution such as e.g. a power law by appropriately choosing model parameters based on some Fourier analysis.

#### 1.3. Related work

Below we point out some connections of our model and results to prior work. More detailed comparisons are provided after we present the main theorems.

The SBM and spectral methods If the node attribute space  $\mathcal{X}$  is a finite set and no edge label is available, then the GSBM reduces to the classical SBM with finite number of blocks. The spectral method and its variants are widely used to recover the underlying clusters under the SBM, see, e.g., McSherry (2001); Coja-Oghlan (2010); Tomozei and Massoulié (2010); Rohe et al. (2011); Chaudhuri et al. (2012). However, the previous analysis relies on the low-rank structure of the edge probability matrix. In contrast, the edge probability matrix under the GSBM is not low-rank, and our analysis is based on establishing a correspondence between the spectrum of a compact operator and the spectrum of a weighted adjacency matrix (see Proposition 4). Similar connection appears before in the context of data clustering considered in von Luxburg et al. (2005), where a graph is constructed based on observed attributes of nodes and clustering based on the graph Laplacian is analyzed. In contrast our setup does not assume the observation of node attributes. Also in our case the observed graphs could be very sparse, while the graphs considered in von Luxburg et al. (2005) are dense.

**Latent space model** If the node attribute space  $\mathcal{X}$  is a finite-dimensional Euclidean space and no edge label is present, then the GSBM reduces to the latent space model, proposed in (Hoff et al. (2002); Handcock et al. (2007)). If we further assume the node attribute space  $\mathcal{X}$  is the probability simplex endowed with Dirichlet distribution with a parameter  $\alpha$ , and B is a bilinear function, then

the SBM reduces to the mixed membership SBM proposed in Airoldi et al. (2008), which is a popular model for studying the overlapping community detection problem.

**Exchangeable random graphs** If we ignore the edge labels, the GSBM fits exactly into the framework of "exchangeable random graphs" and the edge probability function B is known as "graphon" (see e.g., Airoldi et al. (2013) and the references therein). It is pointed out in Bickel and Chen (2009) that some known functions can be used to approximate the graphon, but no analysis is presented. Our spectral algorithm approximates the graphon using the eigenfunctions and the approximation error is determined by the tail of the spectrum of a suitably defined compact operator (see eq. (8)). The exchangeable random graph models with constant average node degrees has been studied in Bollobás et al. (2007), but the focus there is on the phase transition for the emergence of the giant connected component.

**Phase transition if**  $\omega = O(1)$  There is an emerging line of works Decelle et al. (2011); Mossel et al. (2012, 2013); Massoulié (2014); Heimlicher et al. (2012); Lelarge et al. (2013) that try to identify the sharp phase transition threshold for positively correlated clustering in the regime with a bounded average node degree. All these previous works focus on the two communities case. Here we consider the more general case with multiple communities and identify a threshold below which positively correlated clustering is impossible. However, our phase transition threshold is not sharp.

#### 1.4. Notation

For two discrete probability distributions  $\mu$  and  $\nu$  on  $\mathcal{L}$ , let  $\|\mu - \nu\|_{\text{TV}} := \frac{1}{2} \sum_{\ell \in \mathcal{L}} |\mu(\ell) - \nu(\ell)|$ denote the total variation distance. Throughout the paper, we say an event occurs "a.a.s." or "asymptotically almost surely" when it occurs with a probability tending to one as  $n \to \infty$ . We use the standard big O notation. For instance, for two sequences  $\{a_n\}, \{b_n\}, a_n \sim b_n$  means  $\lim_{n\to\infty} \frac{a_n}{b_n} = 1$ .

#### **2.** Spectral reconstruction if $\omega = \Omega(\log n)$

Let  $A \in \{0,1\}^{n \times n}$  denote the adjacency matrix of G and  $L_{ij} \in \mathcal{L}$  denote the label of edge (i, j) in G. Our goal reduces to infer  $\mu_{\sigma_i,\sigma_j}$  based on A and L. In this section, we study a polynomial-time algorithm based on the spectrum of a suitably weighted adjacency matrix. The detailed description is given in Algorithm 1 with four steps.

Step 1 defines the weighted adjacency matrix  $\tilde{A}$  using a random weighing function W of edge labels. Step 2 extracts the top r eigenvalues and eigenvectors of  $\tilde{A}$  for a given integer r. Step 3 embeds n nodes in  $\mathbb{R}^r$  based on the spectrum of  $\tilde{A}$ . Step 4 constructs an estimator of  $\mu_{\sigma_i,\sigma_j}$  using the empirical label distribution on the edges between node j and nodes in the local neighborhood of node i. Note that the random weight function W chosen in Step 1 is the key to exploit the labeling information encoded in G. If  $\nu_{x,y}$  were known, better deterministic weight function could be chosen to allow for sharper estimation, e.g. (Lelarge et al. (2013)). However, no a priori deterministic weight function could ensure consistent estimation irrespective of  $\nu_{x,y}$ . The function  $h_{\epsilon}(x) :=$  $\min(1, \max(0, 2 - x/\epsilon))$  used in Step 4 is a continuous approximation of the indicator function  $\mathbb{I}_{\{x \leq \epsilon\}}$  such that  $h_{\epsilon}(x) = 1$  if  $x \leq \epsilon$  and  $h_{\epsilon} = 0$  if  $x \geq 2\epsilon$ .

Our performance guarantee of Spectral Algorithm 1 is stated in terms of the spectrum of the integral operator defined as

$$Tf(x) := \int_{\mathcal{X}} K(x, y) f(y) P(dy), \tag{5}$$

Algorithm 1 Spectral Algorithm  $(A, L, r, \epsilon)$ 

- 1: (Random Weighing) Let  $W : \mathcal{L} \to [0, 1]$  be a random weighing function, with i.i.d. weights  $W(\ell)$  uniformly distributed on [0, 1]. Define the weighted adjacency matrix as  $\tilde{A}_{ij} = W(L_{ij})$  if  $A_{ij} = 1$ ; otherwise  $\tilde{A}_{ij} = 0$ .
- (Spectral Decomposition) For a given positive integer r, extract the r largest eigenvalues of A
   sorted in decreasing order |λ<sub>1</sub><sup>(n)</sup>| ≥ |λ<sub>2</sub><sup>(n)</sup>| ≥ ··· ≥ |λ<sub>r</sub><sup>(n)</sup>| and the corresponding eigenvectors
   with unit norm v<sub>1</sub>, v<sub>2</sub>,..., v<sub>r</sub> ∈ ℝ<sup>n</sup>.
- 3: (Spectral Embedding) Embed the n nodes in  $\mathbb{R}^r$  by letting

$$z_i := \sqrt{n} \left( \frac{\lambda_1^{(n)}}{\lambda_1^{(n)}} v_1(i), \dots, \frac{\lambda_r^{(n)}}{\lambda_1^{(n)}} v_r(i) \right).$$

$$(2)$$

4: (Label Estimation) For a given small positive parameter  $\epsilon$ , define the estimator  $\hat{\mu}_{ij}$  of  $\mu_{\sigma_i,\sigma_j}$  by letting

$$\hat{\mu}_{ij}(\ell) := \frac{\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) \mathbb{I}_{\{L_{i'j} = \ell\}}}{\epsilon + \sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) A_{i'j}}.$$
(3)

Define the estimator  $\hat{B}_{ij}$  of  $\frac{\omega}{n} B_{\sigma_i,\sigma_j}$  by letting

$$\hat{B}_{ij}(\ell) := \frac{\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) A_{i'j}}{\epsilon + \sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2)}.$$
(4)

where the symmetric kernel K is defined by

$$K(x,y) := \sum_{\ell} W(\ell) \nu_{x,y}(\ell) \in [0, |\mathcal{L}|].$$
(6)

Since K is bounded, the operator T, acting on the function space  $L_2(P)$ , is *compact* and therefore admits a discrete spectrum with finite multiplicity of all of its non-zero eigenvalues (see e.g. Kato (1966) and von Luxburg et al. (2005)). Moreover, any of its eigenfunctions is continuous on  $\mathcal{X}$ . Denote the eigenvalues of operator T sorted in decreasing order by  $|\lambda_1| \ge |\lambda_2| \ge \cdots$  and its corresponding eigenfunctions with unit norm by  $\phi_1, \phi_2, \cdots$ . Define

$$d^{2}(x,x') := \int_{\mathcal{X}} |K(x,y) - K(x',y)|^{2} P(dy).$$
(7)

It is easy to check that with probability 1 with respect to the random choices of  $W(\ell)$ , by the identifiability condition (1), d(x, x') > 0 for all  $x \neq x' \in \mathcal{X}$ . By Minkowski inequality, d(x, x') satisfies the triangle inequality. Therefore, d(x, x') is a distance on  $\mathcal{X}$ . By the definition of  $\lambda_k$  and  $\phi_k$ , we have (the following serie converges in  $L_2(P \times P)$ , see Chapter V.4 in Kato (1966)):

$$K(x,y) = \sum_{k=1}^{\infty} \lambda_k \phi_k(x) \phi_k(y), \tag{8}$$

and thus  $d^{2}(x, x') = \sum_{k=1}^{\infty} \lambda_{k}^{2} (\phi_{k}(x) - \phi_{k}(x'))^{2}$ .

To derive the performance guarantee of Spectral Algorithm 1, we make the following continuity assumption on  $\nu_{x,y}$ . Similar continuity assumptions appeared before in the literature on the latent space model and the exchangeable random graph model (see e.g., (Chatterjee, 2012, Section 4.4) and (Airoldi et al., 2013, Section 2.1)).

**Assumption 1** For every  $\ell \in \mathcal{L}$ ,  $\nu_{x,y}(\ell)$  is continuous on  $\mathcal{X}^2$ , hence by compactness of  $\mathcal{X}$  uniformly continuous. Let  $\psi(\cdot)$  denote a modulus of continuity of all functions  $(x, y) \to \nu_{x,y}(\ell)$  and  $(x, y) \to B_{x,y}$ . That is to say, for all x, x', y, y',

$$|B_{x,y} - B_{x',y'}| \le \psi(d(x,x') + d(y,y'))$$

and similarly for  $\nu_{x,y}(\ell)$ .

Let  $\epsilon_r := \sum_{k>r} \lambda_k^2$  for a fixed integer r, characterizing the tail of the spectrum of the operator T. The following theorem gives an upper bound of the estimation error of  $\hat{\mu}_{ij}$  for most pairs (i, j) in terms of  $\epsilon_r$  and  $\epsilon$ .

**Theorem 2** Suppose Assumption 1 holds. Assume that  $\omega \ge C \log n$  for some universal positive constant C and r chosen in Spectral Algorithm 1 satisfies  $|\lambda_r| > |\lambda_{r+1}|$ . Then a.a.s. the estimators  $\hat{\mu}$  and  $\hat{B}$  given in Spectral Algorithm 1 satisfy

$$B_{\sigma_{i},\sigma_{j}}|\hat{\mu}_{ij}(\ell) - \mu_{\sigma_{i},\sigma_{j}}(\ell)| \leq 2\psi(2|\lambda_{1}|\epsilon) + \frac{1}{|\lambda_{1}|^{2}\epsilon^{2}} \frac{\sqrt{\epsilon_{r}}}{\int_{\mathcal{X}} h_{|\lambda_{1}|\epsilon}(d(\sigma_{i},x))P(dx)} := \eta, \ \forall \ell \in \mathcal{L},$$

$$|\hat{B}_{ij} - \frac{\omega}{n}B_{\sigma_{i},\sigma_{j}}| \leq \frac{\omega}{n}\eta,$$
(9)

for a fraction of at least  $(1 - \sqrt{\epsilon_r})$  of all possible pairs (i, j) of nodes.

Note that if  $\epsilon_r$  goes to 0, the second term in  $\eta$  given by (9) vanishes, and  $\eta$  simplifies to  $2\psi(2|\lambda_1|\epsilon)$ which goes to 0 if  $\epsilon$  further goes to 0. In the case where  $B_{\sigma_i,\sigma_j}$  is strictly positive, Theorem 2 implies that the estimation error of the edge label distribution goes to 0 as successively  $\epsilon_r$  and  $\epsilon$  converge to 0. Note that  $\epsilon$  is a free parameter chosen in Spectral Algorithm 1 and can be made arbitrarily small if  $\epsilon_r$  is sufficiently small. The parameter  $\epsilon_r$  measures how well the compact space  $\mathcal{X}$  endowed with measure P can be approximated by r discrete points, or equivalently how well our general model can be approximated by the labeled stochastic block model with r blocks. The smaller  $\epsilon_r$  is, the more structured, or the more "low-dimensional" our general model is. In this sense, Theorem 2 establishes an interesting connection between the estimation error and the structure present in our general model.

A key part of the proof of Theorem 2 is to show that for any fixed k, the normalized k-th largest eigenvalue  $\lambda_k^{(n)}/\lambda_1^{(n)}$  of the weighted adjacency matrix  $\tilde{A}$  asymptotically converges to  $\lambda_k/\lambda_1$  where  $\lambda_k$  is the k-th eigenvalue of integral operator T, and this is precisely why our spectral embedding given by (2) is defined in a normalized fashion. The following simple example illustrates how we can derive closed form expressions for the spectrum of integral operator T.

**Example 1** Take  $\mathcal{X} = [0,1]$  and P as the Lebesgue measure. Assume unlabeled edges. Let  $B_{x,y} = g(x-y)$  where g is an even (i.e.  $g(-\cdot) = g(\cdot)$ ), 1-periodic function. Denote its Fourier series expansion by  $g(x) = \sum_{k\geq 0} g_k \cos(2\pi kx)$ . For instance, if g(x) = |x| for  $x \in [-1/2, 1/2]$ , then  $g_0 = 1/4$  and  $g_k = [(-1)^k - 1]/(\pi^2 k^2)$  for  $k \geq 1$ . If  $g(x) = \mathbb{I}_{\{-1/4 \leq x \leq 1/4\}}$  for  $x \in [-1/2, 1/2]$ , then  $g_0 = 1/2$  and  $g_k = 2\sin(\pi k/2)/(\pi k)$  for  $k \geq 1$ .

For the above example, Tf = g \* f where \* denotes convolution. Fourier series analysis entails that  $\lambda_k$  must coincide with Fourier coefficient  $g_0$  or  $g_k/2$  for  $k \ge 1$  ( $g_k/2$  appearing twice in the spectrum of T). This example thus gives a general recipe for constructing random graph models with spectrum asymptotic to a pre-specified eigenvalue profile. For g(x) = |x| on [-1/2, 1/2], we find in particular that  $\lambda_1 = 1/4$  and  $|\lambda_{2k}| = |\lambda_{2k+1}| = 1/(\pi^2(2k-1)^2)$ , which is a power-law spectrum with the decaying exponent being 2. For  $g(x) = \mathbb{I}_{\{-1/4 \le x \le 1/4\}}$  on [-1/2, 1/2],  $\lambda_1 = 1/2$ and  $|\lambda_{2k}| = |\lambda_{2k+1}| = 1/(\pi(2k-1))$ , which is a power-law spectrum with the decaying exponent being 1.

**Comparisons to previous work** Theorem 2 provides the first theoretical result on inferring edge label distributions to our knowledge. For estimating edge probabilities, Theorem 2 implies or improves the best known results in several special cases.

For the SBM with finite r blocks,  $\epsilon_r$  is zero. By choosing  $\epsilon$  sufficiently small in Theorem 2, we see that our spectral method is asymptotically correct if  $\omega = \Omega(\log n)$ , which matches with best known bounds (see e.g., Chen and Xu (2014) and the references therein). For the mixed membership SBM with finite r blocks, the best known performance guarantee given by Anandkumar et al. (2013) needs  $\omega$  to be above the order of several log n factors, while Theorem 2 only needs  $\omega$  to be the order of log n. However, Theorem 2 requires the additional spectral gap assumption and needs  $\epsilon_r$  to vanish. Also, notice that Theorem 2 only applies to the setting where the edge probability p within the community exceeds the edge probability q across two different communities by a constant factor, while the best known results in Chen and Xu (2014); Anandkumar et al. (2013) apply to the general setting with any r, p, q.

For the latent space model, Chatterjee (2012) proposed a universal singular value thresholding approach and showed that the edge probabilities can be consistently estimated if  $\omega \ge n^{\frac{k}{k+2}}$  with some Lipschitz condition on B similar to Assumption 1, where k is the dimension of the node attribute space. Our results in Theorem 2 do not depend on the dimension of the node attribute space and only need  $\omega$  to be on the order of  $\log n$ .

For the exchangeable random graph models, a singular value thresholding approach is shown in Chatterjee (2012) to estimate the graphon consistently. More recently, Airoldi et al. (2013) shows that the graphon can be consistently estimated using the empirical frequency of edges in local neighborhoods, which are constructed by thresholding based on the pairwise distances between different rows of the adjacency matrix. All these previous works assume the edge probabilities are constants. In contrast, Theorem 2 applies to much sparser graphs with edge probabilities could be as low as  $\log n/n$ .

#### **3.** Impossibility if $\omega = O(1)$

We have seen in the last section that Spectral Algorithm 1 achieves asymptotically correct inference of edge label distributions so long as  $\omega = \Omega(\log n)$  and  $\epsilon_r = 0$ . In this section, we focus on the sparse regime where  $\omega$  is a constant, i.e., the average node degree is bounded and the number of observed edge labels is only linearly in n. We identify a non-trivial threshold under which it is fundamentally impossible to infer the edge label distributions with an accuracy better than guessing without using the observations.

To derive the impossibility result, let us consider a simple scenario where the compact space  $\mathcal{X} = \{1, \ldots, r\}$  is endowed with a uniform measure  $P, B_{x,y} = \frac{a}{a+b}$  if x = y and  $B_{x,y} = \frac{b}{a+b}$ 

if  $x \neq y$  for two positive constants a, b, and  $\mu_{x,y} = \mu$  if x = y and  $\mu_{x,y} = \nu$  if  $x \neq y$  for two different discrete probability measures  $\mu, \nu$  on  $\mathcal{L}$ . Since  $\omega$  is a constant, the observed labeled graph G is sparse and has a bounded average degree. Similar to the Erdős-Rényi random graph, there are at least  $\Theta(n)$  isolated nodes without neighbors. To infer the edge label distribution between two isolated nodes, the observed labeled graph G does not provide any useful information and thus it is impossible to achieve the asymptotically correct inference of edge label distribution for two isolated nodes. Hence we resort to a less ambitious objective.

**Objective 1** Given any two randomly chosen nodes *i* and *j*, we would like to correctly determine whether the label distribution is  $\mu$  or  $\nu$  with probability strictly larger than 1 - 1/r, which is achievable by always guessing  $\nu$  and is the best one can achieve if no graph information available.

Note that if Objective 1 is not achievable, then the expected estimation error is at least  $\frac{1}{2r} \|\mu - \nu\|_{\text{TV}}$ . One might think that we can always achieve Objective 1 as long as  $\omega > 0$  such that the graph contains a giant connected component, because the labeled graph *G* then could provide extra information. It turns out that this is not the case. Define

$$\tau = \frac{1}{r(a+b)} \sum_{\ell \in \mathcal{L}} |a\mu(\ell) - b\nu(\ell)|.$$
(10)

Let  $\omega_0 = 1/\tau$  and  $\omega_c = \frac{r(a+b)}{a+(r-1)b}$ . Then by definition of  $\tau$ , we have  $\omega_0 > \omega_c$ . Note that when  $\omega > \omega_c$ , the average node degree is larger than one, and thus similar to Erdős-Rényi random graph, G contains a giant connected component. The following theorem shows that Objective 1 is fundamentally impossible if  $\omega < \omega_0$  where  $\omega_0$  is strictly above the threshold  $\omega_c$  for the emergence of the giant connected component.

**Theorem 3** If  $\omega < \omega_0$ , then for any two randomly chosen nodes  $\rho$  and v,

$$\forall x, y \in \{1, \dots, r\}, \ \mathbb{P}(\sigma_{\rho} = x | G, \sigma_v = y) \sim \frac{1}{r} a.a.s$$

The above theorem implies that it is impossible to correctly determine whether two randomly chosen nodes have the same attribute or not with probability larger than 1 - 1/r and thus Objective 1 is fundamentally impossible. In case  $a \neq b$ , it also implies that we cannot correctly determine whether the edge probability between nodes *i* and *j* is  $\frac{a}{a+b}$  or  $\frac{b}{a+b}$  with probability strictly larger than 1-1/r. This indicates the need for a minimum number of observations in order to exploit the information encoded in the labeled graph.

**Comparisons to previous work** To our knowledge, Theorem 2 provides the first impossibility result on inferring edge label distributions and node attributes in the case with multiple communities. The previous work focuses on the case with two community case. If r = 2 and no edge label is available, it is conjectured in Decelle et al. (2011) and later proved in Mossel et al. (2012, 2013); Massoulié (2014) that the positively correlated clustering is feasible if and only if  $(a-b)^2 > 2(a+b)$ , or equivalently,  $\omega > 1/2\tau^2$ . If the edge label is available, it is conjectured in Heimlicher et al. (2012) that the positively correlated clustering is feasible if and only if  $\omega > 1/\tau'$  with

$$\tau' = \frac{1}{2(a+b)} \sum_{\ell \in \mathcal{L}} \frac{(a\mu(\ell) - b\nu(\ell))^2}{a\mu(\ell) + b\nu(\ell)} \le \tau.$$

It is proved in Lelarge et al. (2013) that the positively correlated clustering is infeasible if  $\omega < 1/\tau'$ . Comparing to the previous works, the threshold  $1/\tau$  given by Theorem 3 is not sharp in the special case with two communities.

#### 4. Numerical experiments

In this section, we explore the empirical performance of our Spectral Algorithm 1 based on Example 1. In particular, suppose n = 1500 nodes are uniformly distributed over the space  $\mathcal{X} = [0, 1]$ . Let  $B_{x,y} = g(x - y)$  where g is even, 1-periodic and defined by g(x) = |x| for  $x \in [-1/2, 1/2]$ . Assume unlabeled edges first.

We simulate the spectral embedding given by Step 3 of Algorithm 1 for a fixed observation probability  $\omega/n = 0.6$ . Pick r = 3 in Algorithm 1. Note that the eigenvector  $v_1$  corresponding to the largest eigenvalue is nearly parallel to the all-one vector and thus does not convey any useful information. Therefore, our spectral embedding of n nodes are based on  $v_2$  and  $v_3$ . In particular, let  $z_i = (v_3(i), v_2(i)) \in \mathbb{R}^2$ . As we derived in Section 2, the second and third largest eigenvalues of operator T are given by  $\lambda_2 = \lambda_3 = -1/\pi^2$ , and the corresponding eigenfunctions are given by  $\phi_2(x) = \sqrt{2} \cos(2\pi x)$  and  $\phi_3(x) = \sqrt{2} \sin(2\pi x)$ . Proposition 4 shows that  $z_i$  asymptotically converges to  $f_i = \sqrt{\frac{2}{n}} (\cos(2\pi\sigma_i), \sin(2\pi\sigma_i))$ . We plot  $f_i$  and  $z_i$  in a two-dimensional plane as shown in Fig. 1(a) and Fig. 1(b), respectively. For better illustration, we divide all nodes into ten groups with different colors, where the k-th group consists of nodes with attributes given by  $\frac{1}{n}[100(k-1)+1, 100(k-1)+2, \ldots, 100k]$ . As we can see,  $z_i$  is close to  $f_i$  for most nodes i, which coincides with our theoretical finding.



Figure 1: (a): The spectral embedding given by  $f_i$ ; (b): The spectral embedding given by  $z_i$ .

Then we simulate Spectral Algorithm 1 on estimating the observed edge probability  $\frac{\omega}{n}B_{\sigma_i,\sigma_j}$  between any node pair (i, j) by picking r = 3 and setting  $\epsilon = 0.5$  median  $\{||z_i - z_j||\}$ . We measure the estimation error by the normalized mean square error given by  $||\hat{B} - \frac{\omega}{n}B^*||_F / ||\bar{B} - \frac{\omega}{n}B^*||_F$ , where  $B^*$  is the true edge probability defined by  $B_{ij}^* = B_{\sigma_i,\sigma_j}$ ;  $\hat{B}$  is our estimator defined in (4);  $\bar{B}_{ij}$  is the empirical average edge probability defined by  $\bar{B}_{ij} = \sum_{i'} A_{i,i'} / (n-1)$ . Our simulation result is depicted in Fig. 2(a).



Figure 2: (a): Estimating the observed edge probability; (b): Estimating the edge label distribution.

Next we consider labeled edges with two possible labels +1 or -1 and  $\mu_{x,y}(+1) = 2g(x - y)$ . We simulate Spectral Algorithm 1 for estimating  $\mu_{\sigma_i,\sigma_j}$  between any node pair (i, j) by choosing the weight function as  $W(\pm 1) = \pm 1$ . We again measure the estimation error by the normalized mean square error given by  $\|\hat{\mu} - \frac{\omega}{n}\mu^*\|_F / \|\bar{\mu} - \mu^*\|_F$ , where  $\mu^*$  is the true label distribution defined by  $\mu_{ij}^* = \mu_{\sigma_i,\sigma_j}$ ;  $\hat{\mu}$  is our estimator defined in (3);  $\bar{\mu}_{ij}$  is the empirical label distribution defined by  $\bar{\mu}_{ij}(+1) = \sum_{i'} \mathbb{I}_{\{L_{ii'}=+1\}} / \sum_{i'} A_{ii'}$ . Our simulation result is depicted in Fig. 2(*b*). As we can see from Fig. 2, when  $\omega/n$  is larger than 0.1, our spectral algorithm performs better than the estimator based on the empirical average.

#### 5. Proofs

#### 5.1. Proof of Theorem 2

Our proof is divided into three parts. We first establish the asymptotic correspondence between the spectrum of the weighted adjacency matrix  $\tilde{A}$  and the spectrum of the operator T using Proposition 4. Then, we prove that the estimator of edge label distribution converges to a limit. Finally, we upper bound the total variation distance between the limit and the true label distribution using Proposition 5.

**Proposition 4** Assume that  $\omega \ge C \log n$  for some universal positive constant C and r chosen in Spectral Algorithm 1 satisfies  $|\lambda_r| > |\lambda_{r+1}|$ . Then for k = 1, 2, ..., r+1, almost surely  $\lambda_k^{(n)}/\lambda_1^{(n)} \sim \lambda_k/\lambda_1$ . Moreover, for k = 1, 2, ..., r, almost surely there exist choices of orthonormal eigenfunctions  $\phi_k$  of operator T associated with  $\lambda_k$  such that  $\lim_{n\to\infty} \sum_{i=1}^n (v_k(i) - \frac{1}{\sqrt{n}}\phi_k(\sigma_i))^2 = 0$ .

By Proposition 4, we get the existence of eigenfunctions  $\phi_k$  of T associated with  $\lambda_k$  such that a.a.s., by letting

$$f_m := \left(\frac{\lambda_1}{\lambda_1}\phi_1(\sigma_m), \dots, \frac{\lambda_r}{\lambda_1}\phi_r(\sigma_m)\right),$$

we have

$$\sum_{n=1}^{n} ||z_m - f_m||_2^2 = \sum_{m=1}^{n} \sum_{k=1}^{r} \left( \sqrt{n} \frac{\lambda_k^{(n)}}{\lambda_1^{(n)}} v_k(m) - \frac{\lambda_k}{\lambda_1} \phi_k(\sigma_m) \right)^2 = o(n).$$

By Markov's inequality,

$$\frac{1}{n} |\{m \in \{1, \dots, n\} : ||z_m - f_m||_2 \ge \delta_n\}| \le \frac{\sum_{m=1}^n ||z_m - f_m||_2^2}{n\delta_n^2} = \frac{1}{\delta_n^2} o(1).$$

Note that  $\delta_n$  can be chosen to decay to zero with n sufficiently slowly so that the right-hand side of the above is o(1). We call nodes m satisfying  $||z_m - f_m||_2 \ge \delta_n$  "bad nodes". Let  $\mathcal{I}$  denote the set of "bad nodes". It follows from the last display that  $|\mathcal{I}| = o(n)$ . Let  $\mathcal{J}$  denote the set of nodes with at least  $\gamma_n$  fraction of edges directed towards "bad nodes", i.e.,

$$\mathcal{J} = \{j : |\{i \in \mathcal{I} : A_{ij} = 1\}| \ge \gamma_n |\{i : A_{ij} = 1\}|\}$$

Note that the average node degree in G is  $\Theta(\omega)$ . Since  $\omega \ge C \log n$  by assumption, it follows from the Chernoff bound that the observed node degree is  $\Theta(\omega)$  with high probability. Therefore, we can choose  $\gamma_n$  decaying to zero while still having  $|\mathcal{J}| = o(n)$ , i.e., all but a vanishing fraction of nodes have at most  $\gamma_n$  fraction of edges directed towards "bad nodes". We have thus performed an embedding of n nodes in  $\mathbb{R}^r$  such that for  $m, m' \notin \mathcal{I}$ ,

$$||z_m - z_{m'}||_2 = \frac{1}{|\lambda_1|} d_r(\sigma_m, \sigma_{m'}) + O(\delta_n),$$
(11)

where pseudo-distance  $d_r$  is defined by  $d_r^2(x, x') := \sum_{k=1}^r \lambda_k^2 \left( \phi_k(x) - \phi_k(x') \right)^2$ .

The remainder of the proof exploits this embedding and the fact that pseudo-distance  $d_r$  and distance d are apart by at most  $\epsilon_r$  in some suitable sense. For a randomly selected pair of nodes (i, j), one has a.a.s.  $i \notin \mathcal{I}$  and  $j \notin \mathcal{J}$ . Therefore, node j has at most  $\gamma_n = o(1)$  fraction of edges directed towards "bad nodes". Hence, by (11),

$$\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) \mathbb{I}_{\{L_{i'j} = \ell\}} = \sum_{i'} \mathbb{I}_{\{L_{i'j} = \ell\}} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, \sigma_{i'}) + O(\delta_n) \right) + O(\omega\gamma_n), \quad (12)$$

and

$$\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) A_{i'j} = \sum_{i'} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, \sigma_{i'}) + O(\delta_n) \right) A_{i'j} + O(\omega\gamma_n).$$
(13)

The first term in the R.H.S. of (12) is a sum of i.i.d. bounded random variables with mean given by

$$\frac{\omega}{n} \int_{\mathcal{X}} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, x) + O(\delta_n) \right) \nu_{x,\sigma_j}(\ell) P(dx).$$

Since  $\omega \ge C \log n$  by assumption, it follows from the Bernstein inequality that a.a.s.

$$\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) \mathbb{I}_{\left\{L_{i'j} = \ell\right\}} = (1 + o(1))\omega \int_{\mathcal{X}} h_{|\lambda_1|\epsilon} \left(d_r(\sigma_i, x) + O(\delta_n)\right) \nu_{x,\sigma_j}(\ell) P(dx) + O(\omega\gamma_n).$$

$$(14)$$

The first term in the R.H.S. of (13) is a sum of i.i.d. bounded random variables with mean given by

$$\frac{\omega}{n} \int_{\mathcal{X}} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, x) + O(\delta_n) \right) B_{x, \sigma_j} P(dx).$$

It again follows from the Bernstein inequality that a.a.s.

$$\sum_{i'} h_{\epsilon}(||z_{i'} - z_i||_2) A_{i'j} = (1 + o(1))\omega \int_{\mathcal{X}} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, x) + O(\delta_n) \right) B_{x,\sigma_j} P(dx) + O(\omega\gamma_n).$$
(15)

Note that  $h_{\epsilon}(x)$  is a continuous function in x. Therefore,

$$\lim_{n \to \infty} h_{|\lambda_1|\epsilon} \left( d_r(\sigma_i, x) + O(\delta_n) \right) = h_{|\lambda_1|\epsilon} (d_r(\sigma_i, x)).$$

By the dominated convergence theorem, it follows from (3), (14), (15) that a.a.s.

$$\hat{\mu}_{i,j}(\ell) \sim \frac{\int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d_r(\sigma_i, x))\nu_{x,\sigma_j}(\ell)P(dx)}{\int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d_r(\sigma_i, x))B_{x,\sigma_j}P(dx)} := \mu_{i,j}^*(\ell).$$
(16)

Similarly, we have a.a.s.

$$\hat{B}_{i,j}(\ell) \sim \frac{\omega}{n} \frac{\int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d_r(\sigma_i, x)) B_{x,\sigma_j} P(dx)}{\int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d_r(\sigma_i, x)) P(dx)} := B_{i,j}^*.$$
(17)

The following proposition upper bounds the difference between the limit  $\mu_{i,j}^*(\ell)$  (resp.  $B_{i,j}^*(\ell)$ ) and  $\mu_{\sigma_i,\sigma_j}(\ell)$  (resp.  $B_{\sigma_i,\sigma_j}$ ).

**Proposition 5** Suppose Assumption 1 holds. Then there exists a fraction of at least  $(1 - \sqrt{\epsilon_r})$  of all possible pairs (i, j) of nodes such that

$$B_{\sigma_{i},\sigma_{j}}|\mu_{i,j}^{*}(\ell) - \mu_{\sigma_{i},\sigma_{j}}(\ell)| \leq 2\psi(2|\lambda_{1}|\epsilon) + \frac{1}{|\lambda_{1}|^{2}\epsilon^{2}} \frac{\sqrt{\epsilon_{r}}}{\int_{\mathcal{X}} h_{|\lambda_{1}|\epsilon}(d(\sigma_{i},x))P(dx)} := \eta, \ \forall \ell \in \mathcal{L},$$

$$|B_{i,j}^{*} - \frac{\omega}{n}B_{\sigma_{i},\sigma_{j}}| \leq \frac{\omega}{n}\eta.$$
(18)

Applying Proposition 5, our theorem then follows.

#### 5.2. Proof of Theorem 3

Proof of Theorem 3 relies on a nice coupling between the local neighborhood of  $\rho$  with a simple labeled Galton-Watson tree. It is well-known that the local neighborhood of a node in the sparse graph is "tree-like". In the case with r = 2, the coupling result is first studied in Mossel et al. (2012) and generalized to the labeled tree in Lelarge et al. (2013). In this paper, we extend the coupling result to any finite  $r \ge 2$ .

Let  $d = \omega \frac{a+(r-1)b}{r(a+b)}$  and consider a labeled Galton-Watson tree  $\mathcal{T}$  with Poisson offspring distribution with mean d. The attribute of root  $\rho$  is chosen uniformly at random from  $\mathcal{X}$ . For each child node, independently of everything else, it has the same attribute with its parent with probability  $\frac{a}{a+(r-1)b}$  and one of r-1 different attributes with probability  $\frac{b}{a+(r-1)b}$ . Every edge between the child and its parent is independently labeled with distribution  $\mu$  if they have the same attribute and with distribution  $\nu$  otherwise.

The labeled Galton-Watson tree  $\mathcal{T}$  can also be equivalently described as follows. Each edge is independently labeled at random according to the probability distribution  $\mathbb{P}(\ell) = \frac{a\mu(\ell) + (r-1)b\nu(\ell)}{a + (r-1)b}$ .

The attribute of root  $\rho$  is first chosen uniformly at random from  $\mathcal{X}$ . Then, for each child node, independently of everything else, it has the same attribute with its parent with probability  $1 - (r - 1)\epsilon(\ell)$  and one of r - 1 different attributes with probability  $\epsilon(\ell)$ , where

$$\epsilon(\ell) = \frac{b\nu(\ell)}{a\mu(\ell) + (r-1)b\nu(\ell)}.$$
(19)

Recall that  $G_R$  denote the neighborhood of  $\rho$  in G within distance R and  $\partial G_R$  denote the nodes at the boundary of  $G_R$ . Let  $\mathcal{T}_R$  denote the tree  $\mathcal{T}$  up to depth R and  $\partial \mathcal{T}_R$  denote the set of leaf nodes of  $\mathcal{T}_R$ . The following lemma similar to coupling lemmas in Mossel et al. (2012) and Lelarge et al. (2013) shows that  $G_R$  can be coupled with the labeled Galton-Watson tree  $\mathcal{T}_R$ .

**Lemma 6** Let  $R = \theta \log n$  for some small enough constant  $\theta > 0$ , then there exists a coupling such that a.a.s.  $(G_R, \sigma_{G_R}) = (\mathcal{T}_R, \sigma_{\mathcal{T}_R})$ , where  $\sigma_{G_R}$  denote the node attributes on the subgraph  $G_R$ .

For the labeled Galton-Watson tree, we show that if  $\omega < \omega_0$ , then the attributes of leaf nodes are asymptotically independent with the attribute of root.

**Lemma 7** Consider a labeled Galton-Waltson tree  $\mathcal{T}$  with  $\omega < \omega_0$ . Then as  $R \to \infty$ ,

$$\forall x \in \{1, \dots, r\}, \ \mathbb{P}(\sigma_{\rho} = x | \mathcal{T}, \sigma_{\partial T_R}) \to \frac{1}{r} a.a.s$$

By exploiting Lemma 6 and Lemma 7, we give our proof of Theorem 3. By symmetry,  $\mathbb{P}[\sigma_{\rho} = x|G, \sigma_{v} = y] = \mathbb{P}[\sigma_{\rho} = x'|G, \sigma_{v} = y]$  for  $x, x' \neq y$  and  $x \neq x'$ . Therefore, we only need to show that  $\mathbb{P}[\sigma_{\rho} = y|G, \sigma_{v} = y] \sim 1/r$  for any  $y \in \mathcal{X}$  and it further reduces to showing that

$$\mathbb{P}[\sigma_{\rho} = y | G, \sigma_{v} = y, \sigma_{\partial G_{R}}] \sim 1/r.$$
(20)

Let  $R = \theta \log n$  be as in Lemma 6 such that  $G_R = o(\sqrt{n})$  and thus  $v \notin G_R$  a.a.s.. Lemma 4.7 in Mossel et al. (2012) shows that  $\sigma_{\rho}$  is asymptotically independent with  $\sigma_v$  conditional on  $\sigma_{\partial G_R}$ . Hence,  $\mathbb{P}[\sigma_{\rho} = y|G, \sigma_v = y, \sigma_{\partial G_R}] \sim \mathbb{P}[\sigma_{\rho} = y|G, \sigma_{\partial G_R}]$ . Also, note that  $\mathbb{P}(\sigma_{\rho} = y|G, \sigma_{\partial G_R}) = \mathbb{P}(\sigma_{\rho} = y|G_R, \sigma_{\partial G_R})$ . Lemma 6 implies that  $\mathbb{P}(\sigma_{\rho} = y|G_R, \sigma_{G_R}) \sim \mathbb{P}(\sigma_{\rho} = y|T_R, \sigma_{\partial T_R})$ , and by Lemma 7,  $\mathbb{P}(\sigma_{\rho} = y|T_R, \sigma_{\partial T_R}) \sim \frac{1}{r}$ . Therefore, equation (20) holds.

#### Acknowledgments

M.L acknowledges the support of the French Agence Nationale de la Recherche (ANR) under reference ANR-11-JS02-005-01 (GAP project). J. X. acknowledges the support of NSF ECCS 10-28464.

#### References

- Edoardo M. Airoldi, David M. Blei, Stephen E. Fienberg, and Eric P. Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9:1981–2014, 2008.
- Edoardo M Airoldi, Thiago B Costa, and Stanley H Chan. Stochastic blockmodel approximation of a graphon: Theory and consistent estimation. In *Advances in Neural Information Processing Systems 26*, pages 692–700, 2013.

- Animashree Anandkumar, Rong Ge, Daniel Hsu, and Sham Kakade. A tensor spectral approach to learning mixed membership community models. In *COLT*, pages 867–881, 2013.
- Peter J. Bickel and Aiyou Chen. A nonparametric view of network models and Newman-Girvan and other modularities. *Proceedings of the National Academy of Sciences*, 2009.
- Béla Bollobás, Svante Janson, and Oliver Riordan. The phase transition in inhomogeneous random graphs. *Random Struct. Algorithms*, 31(1):3–122, August 2007.
- Sourav Chatterjee. Matrix estimation by universal singular value thresholding. *arxiv:1212.1247*, 2012.
- Kamalika Chaudhuri, Fan Chung Graham, and Alexander Tsiatas. Spectral clustering of graphs with general degrees in the extended planted partition model. *Journal of Machine Learning Research*, 23:35.1–35.23, 2012.
- Yudong Chen and Jiaming Xu. Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices. *arxiv*:1402.1267, 2014.
- Amin Coja-Oghlan. Graph partitioning via adaptive spectral techniques. Comb. Probab. Comput., 19(2):227–284, 2010.
- Chandler Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. SIAM Journal on Numerical Analysis, 7(1):pp. 1–46, 1970.
- Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications. *Physics Review E*, 84:066106, 2011.
- William Evans, Claire Kenyon, Yuval Peres, and Leonard J. Schulman. Broadcasting on trees and the Ising model. *The Annals of Applied Probability*, 10(2):410–433, 2000.
- Uriel Feige and Eran Ofek. Spectral techniques applied to sparse random graphs. *Random Struct. Algorithms*, 27(2):251–275, Sept. 2005.
- Santo Fortunato. Community detection in graphs. arXiv:0906.0612, 2010.
- Mark S. Handcock, Adrian E. Raftery, and Jeremy M. Tantrum. Model-based clustering for social networks. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 170(2):301–354, 2007.
- Simon Heimlicher, Marc Lelarge, and Laurent Massoulié. Community detection in the labelled stochastic block model. *arXiv:1209.2910*, 2012.
- Peter D. Hoff, Adrian E. Raftery, and Mark S. Handcock. Latent space approaches to social network analysis. *Journal of the American Statistical Association*, 97:1090+, December 2002.
- Tosio Kato. Perturbation Theory for Linear Operators. Springer, Berlin, 1966.
- Vladimir I. Koltchinskii. Asymptotics of spectral projections of some random matrices approximating integral operators. *Progress in Probability*, 1998.

- Marc Lelarge, Laurent Massoulié, and Jiaming Xu. Reconstruction in the labeled stochastic block model. In *Information Theory Workshop*, Sept. 2013.
- Laurent Massoulié. Community detection thresholds and the weak Ramanujan property. In *STOC* 2014: 46th Annual Symposium on the Theory of Computing, pages 1–10, United States, 2014.
- Frank McSherry. Spectral partitioning of random graphs. In 42nd IEEE Symposium on Foundations of Computer Science, pages 529 537, Oct. 2001.
- Elchanan Mossel. Survey information flows on trees. *DIMACS series in discrete mathematics and theoretical computer science*, pages 155–170, 2004.
- Elchanan Mossel, Joe Neeman, and Allan Sly. Stochastic block models and reconstruction. *arXiv:1202.1499*, 2012.
- Elchanan Mossel, Joe Neeman, and Allan Sly. A proof of the block model threshold conjecture. *arxiv*:1311.4115, 2013.
- Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *The Annals of Statistics*, 39(4):1878–1915, 2011.
- Dan-Cristian Tomozei and Laurent Massoulié. Distributed user profiling via spectral methods. SIG-METRICS Perform. Eval. Rev., 38(1):383–384, June 2010.
- Ulrike von Luxburg, Olivier Bousquet, and Mikhail Belkin. On the convergence of spectral clustering on random samples: the normalized case. *NIPS*, 2005.

#### Appendix A. Proof of Proposition 4

We first introduce notations used in the proof. Several norms on matrices will be used. The spectral norm of a matrix X is denoted by ||X|| and equals the largest singular value. The Frobenius norm of a matrix X is denoted by  $||X||_F$  and equals the square root of the sum of squared singular values. It follows that  $||X||_F \leq \sqrt{r} ||X||_2$  if X is of rank r. For vectors, the only norm that will be used is the usual  $l_2$  norm, denoted as  $||x||_2$ . Introduce a  $n \times n$  matrix  $\hat{A}$  defined by  $\hat{A}_{ij} = K(\sigma_i, \sigma_j)$ . Recall that r is a fixed positive integer in Spectral Algorithm 1. Denote r the largest eigenvalues of  $\hat{A}$  sorted in decreasing order by  $|\lambda_1'^{(n)}| \geq \cdots \geq |\lambda_r'^{(n)}|$  of  $\hat{A}$ . Let  $v'_1, \ldots, v'_r \in \mathbb{R}^n$  be the corresponding eigenvectors with unit norm. An overview of the proof is shown in Fig. 3.

$$\begin{array}{c} \tilde{A} \xleftarrow{\text{Lemma 9 and 10}} \tilde{A} \xleftarrow{\text{Lemma 8}} T \\ \{\lambda_k^{(n)}, v_k\}_{k=1}^r \quad \{\lambda_k'^{(n)}, v_k'\}_{k=1}^r \quad \{\lambda_k, \phi_k\}_{k=1}^r \end{array}$$

Figure 3: Proof outline for showing the asymptotic correspondence between the spectrum of  $\hat{A}$  and that of T.

Lemma 8 follows from the results of Koltchinskii (1998) and their application as per Theorem 4 and Theorem 5 of von Luxburg et al. (2005).

**Lemma 8** Under our assumptions on operator T, for k = 1, 2, ..., r + 1, almost surely  $\frac{1}{n}\lambda_k^{\prime(n)} \sim \lambda_k$ , and there exist choices of orthonormal eigenfunctions  $\phi_k$  of operator T associated with  $\lambda_k$  such that  $\lim_{n\to\infty} \sum_{i=1}^n (v'_k(i) - \frac{1}{\sqrt{n}}\phi_k(\sigma_i))^2 = 0$ .

Lemma 9 gives sharp controls for the spectral norm of random symmetric matrices with bounded entries initially developed by Feige and Ofek (2005) and extended by Tomozei and Massoulié (2010) and Chatterjee (2012).

**Lemma 9** Let M be a random symmetric matrix with entries  $M_{ij}$  independent up to symmetry,  $M_{ij} \in [0,1]$  and such that  $\mathbb{E}[M_{ij}] = \omega/n$ . If  $\omega \ge C \log n/n$  for a universal positive constant C, then for all c > 0 there exists c' > 0 such that with probability at least  $1 - n^{-c}$ , one has

$$\|M - \mathbb{E}[M]\| \le c'\sqrt{\omega}.\tag{21}$$

Lemma 10, a consequence of the famous Davis-Kahan  $\sin \theta$  Theorem Davis and Kahan (1970), controls the perturbation of eigenvectors of perturbed matrices.

**Lemma 10** For two symmetric matrices M, M' and orthonormal eigenvectors  $(u_1, \ldots, u_r)$  (respectively  $u'_1, \ldots, u'_r$ ) associated with the r leading eigenvalues of M (respectively M'), denoting  $U = [u_1, \ldots, u_r], U' = [u'_1, \ldots, u'_r]$ , there exists an orthogonal  $r \times r$  matrix O such that

$$||U - U'O|| \le \frac{\sqrt{2}||M - M'||}{|\theta_r| - |\theta_{r+1}| - ||M - M'||},$$
(22)

where  $\theta_k$  is the k-th largest eigenvalue of M' in absolute value.

We omit proofs of lemmas which can be found in the mentioned literature. Next, we present the proof of our proposition. Applying Lemma 10 to  $M = \tilde{A}$  and  $M' = (\omega/n)\hat{A}$ , then we have  $U = [v_1, \ldots, v_r], U' = [v'_1, \ldots, v'_r]$ , and  $\theta_k = (\omega/n)\lambda'^{(n)}_k$  for  $k = 1, \ldots, r + 1$ . By Lemma 9 and observing that  $\mathbb{E}[\tilde{A}] = (\omega/n)\hat{A}$ , it readily follows that  $||M - M'|| = O(\sqrt{\omega})$  with high probability. By Weyl's inequality, we have  $|\lambda_k^{(n)} - \theta_k| \leq ||M - M'|| = O(\sqrt{\omega})$ . Moreover, by Lemma 8, for  $k = 1, \ldots, r+1, \theta_k \sim \omega\lambda_k$ . Hence,  $\lambda_k^{(n)}/\lambda_1^{(n)} = \lambda_k/\lambda_1 + O(1/\sqrt{\omega})$ . By assumption,  $|\lambda_r| > |\lambda_{r+1}|$ , and thus the right-hand side of (22) is  $O(1/\sqrt{\omega})$ . Note that U, U'O are of rank r, it follows that

$$||U - U'O||_F \le \sqrt{2r} ||U - U'O|| = O(1/\sqrt{\omega}).$$

Therefore, by Lemma 8, there exist choices of orthonormal eigenfunctions  $\phi_k$  of operator T associated with  $\lambda_k$  such that  $\lim_{n\to\infty} \sum_{i=1}^n (v_k(i) - \frac{1}{\sqrt{n}}\phi_k(\sigma_i))^2 = 0$  for  $k = 1, \ldots, r$ .

#### **Appendix B.** Proof of Proposition 5

The main idea of proof is to show that the pseudo-distance  $d_r$  is close to distance d in an appropriate sense. By definition,  $d(x, x') \ge d_r(x, x')$  and moreover,

$$\int_{\mathcal{X}^2} [d^2(x, x') - d_r^2(x, x')] P(dx) P(dx') = \sum_{k>r} \lambda_k^2 \int_{\mathcal{X}^2} \left(\phi_k(x) - \phi_k(x')\right)^2 \le 2 \sum_{k>r} \lambda_k^2 = 2\epsilon_r.$$

Define  $d_{>r}^2(x, x') = d^2(x, x') - d_r^2(x, x')$ . Markov's inequality entails that

$$\int_{\mathcal{X}} P(dx') \mathbb{I}_{\left\{\int_{\mathcal{X}} d_{>r}^2(x,x') P(dx) \ge 2\sqrt{\epsilon_r}\right\}} \le \sqrt{\epsilon_r}$$

Note that the following inequalities hold

$$0 \le h_{\epsilon}(d_r(x, x')) - h_{\epsilon}(d(x, x')) \le \frac{1}{2\epsilon^2} \left[ d^2(x, x') - d_r^2(x, x') \right] = \frac{d_{>r}^2(x, x')}{2\epsilon^2}.$$
 (23)

By the previous application of Markov's inequality, for a fraction of at least  $1 - \sqrt{\epsilon_r}$  of nodes *i*, it holds that

$$\int_{\mathcal{X}} d_{>r}^2(x,\sigma_i) P(dx) \le 2\sqrt{\epsilon_r}.$$

Combined with the previous Lipschitz property (23) and the definition of  $\mu^*$  given by (16), this entails that for a fraction of at least  $1 - \sqrt{\epsilon_r}$  nodes *i*, one has

$$\frac{a}{b+\frac{\sqrt{\epsilon_r}}{|\lambda_1|^2\epsilon^2}} \leq \mu_{i,j}^*(\ell) \leq \frac{a+\frac{\sqrt{\epsilon_r}}{|\lambda_1|^2\epsilon^2}}{b},$$

where we have introduced the following notations

$$a = \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(x,\sigma_i)\nu_{x,\sigma_j}(\ell)P(dx), \ b = \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(x,\sigma_i)B_{x,\sigma_j}P(dx))$$

Define

$$a' = \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(\sigma_i, x))\nu_{\sigma_i, \sigma_j}(\ell)P(dx), \ b' = \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(\sigma_i, x))B_{\sigma_i, \sigma_j}P(dx)$$

Then,  $\mu_{\sigma_i,\sigma_j}(\ell) = a'/b'$ . Note that for positive constants  $c_1 \leq c_2, c_3 \leq c_4$ ,  $|\frac{c_1}{c_2} - \frac{c_3}{c_4}| \leq \frac{1}{c_4}(|c_1 - c_3| + |c_2 - c_4|)$ . Hence,

$$|\mu_{i,j}^*(\ell) - \mu_{\sigma_i,\sigma_j}(\ell)| \le \frac{|a-a'| + |b-b'| + \frac{\sqrt{\epsilon_r}}{|\lambda_1|^2 \epsilon^2}}{b'}$$

By assumption 1, for all x, x', y, y',

$$|B_{x,y} - B_{x',y'}| \le \psi(d(x,x') + d(y,y'))$$

and similarly for  $\nu_{x,y}(\ell)$ . Therefore,

$$\begin{aligned} |a-a'|+|b-b'| &\leq \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(\sigma_i, x)) \left[ |\nu_{x,\sigma_i}(\ell) - \nu_{\sigma_i,\sigma_j}(\ell)| + |B_{x,\sigma_i} - B_{\sigma_j,\sigma_i}| \right] P(dx) \\ &\leq 2\psi(2|\lambda_1|\epsilon) \int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(\sigma_i, x)) P(dx). \end{aligned}$$

It follows that

$$B_{\sigma_i,\sigma_j}|\mu_{i,j}^*(\ell) - \mu_{\sigma_i,\sigma_j}(\ell)| \le 2\psi(2|\lambda_1|\epsilon) + \frac{\sqrt{\epsilon_r}}{|\lambda_1|^2\epsilon^2} \frac{1}{\int_{\mathcal{X}} h_{|\lambda_1|\epsilon}(d(\sigma_i,x))P(dx)} = \eta.$$

The right-hand side goes to zero as one lets successively  $\epsilon_r$  then  $\epsilon$  go to zero. Similarly, we can show  $|B_{i,j}^* - \frac{\omega}{n}B_{\sigma_i,\sigma_j}| \leq \frac{\omega}{n}\eta$ .

#### Appendix C. Proof of Lemma 7

The proof technique is adapted from Section 5 in Mossel (2004). Consider two distributions on the labeled Galton-Watson tree, one with the attribute of the root being x, and one with the attribute of the root being  $y \neq x$ . We can couple the two distributions in the following way: if the two distributions agree on the attribute of node v, then couple them together such that they also agree for all the children of v; if the two distributions do not agree on the attribute of node v, use the optimal coupling to make them agree as much as possible for each children of v. For each children w with  $L_{vw} = \ell$ , it is easy to check that under the optimal coupling, the two distributions will not agree on the attribute of w with probability  $|1 - r\epsilon(\ell)|$ , where  $\epsilon(\ell)$  is defined in (19). Hence, the non-coupled nodes grow as a branching process with the branching number given by

$$\omega \sum_{\ell} \frac{a\mu(\ell) + (r-1)b\nu(\ell)}{r(a+b)} |1 - r\epsilon(\ell)| = \omega\tau.$$

It is well known that if the branching number  $\omega \tau < 1$ , the branching process will eventually die a.a.s. Thus as  $R \to \infty$ , a.a.s.

$$\mathbb{P}(\sigma_{\partial T_R} | \mathcal{T}, \sigma_{\rho} = x) = \mathbb{P}(\sigma_{\partial T_R} | \mathcal{T}, \sigma_{\rho} = y).$$

By Bayes' formula, the theorem follows.

#### **Appendix D. Bernstein Inequality**

**Theorem 11** Let  $X_1, \ldots, X_n$  be independent random variables such that  $|X_i| \leq M$  almost surely. Let  $\sigma_i^2 = Var(X_i)$  and  $\sigma^2 = \sum_{i=1}^n \sigma_i^2$ , then

$$\mathbb{P}(\sum_{i=1}^{n} X_i \ge t) \le \exp\left(\frac{-t^2}{2\sigma^2 + \frac{2}{3}Mt}\right).$$

It follows then

$$\mathbb{P}(\sum_{i=1}^n X_i \ge \sqrt{2\sigma^2 u} + \frac{2Mu}{3}) \le e^{-u}.$$