

# Impact of regularization on Spectral Clustering

Antony Joseph\* and Bin Yu†

July 22, 2014

## Abstract

The performance of spectral clustering can be considerably improved via regularization, as demonstrated empirically in Amini et al. [2]. Here, we provide an attempt at quantifying this improvement through theoretical analysis. Under the stochastic block model (SBM), and its extensions, previous results on spectral clustering relied on the minimum degree of the graph being sufficiently large for its good performance. By examining the scenario where the regularization parameter  $\tau$  is large we show that the minimum degree assumption can potentially be removed. As a special case, for an SBM with two blocks, the results require the maximum degree to be large (grow faster than  $\log n$ ) as opposed to the minimum degree. More importantly, we show the usefulness of regularization in situations where not all nodes belong to well-defined clusters. Our results rely on a ‘bias-variance’-like trade-off that arises from understanding the concentration of the sample Laplacian and the eigen gap as a function of the regularization parameter. As a byproduct of our bounds, we propose a data-driven technique *DKest* (standing for estimated Davis-Kahan bounds) for choosing the regularization parameter. This technique is shown to work well through simulations and on a real data set.

## 1 Introduction

The problem of identifying communities (or clusters) in large networks is an important contemporary problem in statistics. Spectral clustering is one of the more popular techniques for such a purpose, chiefly due to its computational advantage and generality of application. The algorithm’s generality arises from the fact that it is not tied to any modeling assumptions on the data, but is rooted in intuitive measures of community structure such as *sparsest cut* based measures [11], [24], [16], [20]. Other examples of applications of spectral clustering include manifold learning [4], image segmentation [24], and text mining [9].

---

\*Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Statistics, University of California, Berkeley. email: AntonyJoseph@lbl.gov

†Department of Statistics and EECS, University of California, Berkeley. email: binyu@stat.berkeley.edu

The canonical nature of spectral clustering also generates interest in variants of the technique. Here, we attempt to better understand the impact of regularized forms of spectral clustering for community detection in networks. In particular, we focus on the regularized spectral clustering (RSC) procedure proposed in Amini et al. [2]. Their empirical findings demonstrates that the performance of the RSC algorithm, in terms of obtaining the correct clusters, is significantly better for certain values of the regularization parameter. An alternative form of regularization was studied in Chaudhuri et al. [7] and Qin and Rohe [22].

This paper provides an attempt to provide a theoretical understanding for the regularization in the RSC algorithm. We also propose a practical scheme for choosing the regularization parameter based on our theoretical results. Our analysis focuses on the Stochastic Block Model (SBM) and an extension of this model. Below are the three main contributions of the paper.

- (a) We attempt to understand regularization for the stochastic block model. In particular, for a graph with  $n$  nodes, previous theoretical analyses for spectral clustering, under the SBM and its extensions, [23],[7], [25], [10] assumed that the minimum degree of the graph scales at least by a polynomial power of  $\log n$ . Even when this assumption is satisfied, the dependence on the minimum degree is highly restrictive when it comes to making inferences about cluster recovery. Our analysis provides cluster recovery results that potentially do not depend on the above mentioned constraint on the minimum degree. As an example, for an SBM with two blocks (clusters), our results require that the maximum degree be large (grow faster than  $\log n$ ) rather than the minimum degree. This is done in Section 3.
- (b) We demonstrate that regularization has the potential of addressing a situation where the lower degree nodes do not belong to well-defined clusters. Our results demonstrate that choosing a large regularization parameter has the effect of removing these relatively lower degree nodes. Without regularization, these nodes would hamper with the clustering of the remaining nodes in the following way: In order for spectral clustering to work, the top eigenvectors - that is, the eigenvectors corresponding to the largest eigenvalues of the Laplacian - need to be able to discriminate between the clusters. Due to the effect of nodes that do not belong to well-defined clusters these top eigenvectors do not necessarily discriminate between the clusters with ordinary spectral clustering. This is done in Section 4
- (c) Although our theoretical results deal with the ‘large’  $\tau$  case, it is observed empirically that moderate values of  $\tau$  may produce better clustering performance. Consequently, in Section 5 we propose *DKest*, a data dependent procedure for choosing the regularization parameter. We demonstrate that this works well through simulations and on a real data set. This is in Section 5.

Our theoretical results involve understanding the trade-offs between the *eigen gap* and the concentration of the sample Laplacian when viewed as a

function of the regularization parameter. Assuming that there are  $K$  clusters, the eigen gap refers to the gap between the  $K$ -th smallest eigenvalue and the remaining eigenvalues. An adequate gap ensures that the sample eigenvectors can be estimated well ([26], [20], [16]) which leads to good cluster recovery. The adequacy of an eigen gap for cluster recovery is in turn determined by the concentration of the sample Laplacian.

In particular, a consequence of the Davis-Kahan theorem [5] is that if the spectral norm of the difference of the sample and population Laplacians is small compared to the eigen gap then the top  $K$  eigenvector can be estimated well. Denoting  $\tau$  as the regularization parameter, previous theoretical analyses of regularization ([7], [23]) provided high-probability bounds on this spectral norm. These bounds have a  $1/\sqrt{\tau}$  dependence on  $\tau$ , for large  $\tau$ . In contrast, our high probability bounds behave like  $1/\tau$ , for large  $\tau$ . We also demonstrate that the eigen gap behaves like  $1/\tau$  for large  $\tau$ . The end result is that we show that one can get a good understanding of the impact of regularization by understanding the situation where  $\tau$  goes to infinity. This also explains empirical observations in [2], [22] where it was seen that performance of regularized spectral clustering does not change for  $\tau$  beyond a certain value. Our procedure for choosing the regularization parameter works by providing estimates of the Davis-Kahan bounds over a grid of values of  $\tau$  and then choosing the  $\tau$  that minimizes these estimates.

The paper is divided as follows. In the next subsection we discuss preliminaries. In particular, in Subsection 1.1 we review the RSC algorithm of [2], and also discuss the other forms of regularization in literature. In Section 2 we review the stochastic block model. Our theoretical results, described in (a) and (b) above, are provided in Sections 3 and 4. Section 5 describes our *DKest* data dependent method for choosing the regularization parameter.

## 1.1 Regularized spectral clustering

In this section we review the regularized spectral clustering (RSC) algorithm of Amini et al. [2].

We first introduce some basic notation. A graph with  $n$  nodes and edge set  $E$  is represented by the  $n \times n$  symmetric adjacency matrix  $A = ((A_{ij}))$ , where  $A_{ij} = 1$  if there is an edge between  $i$  and  $j$ , otherwise  $A_{ij}$  is 0. In other words, for  $1 \leq i, j \leq n$ ,

$$A_{ij} = \begin{cases} 1, & \text{if } (i, j) \in E \\ 0, & \text{otherwise} \end{cases}.$$

Given such a graph, the typical community detection problem is synonymous with finding a partition of the nodes. A good partitioning would be one in which there are fewer edges between the various components of the partition, compared to the number of edges within the components. Various measures for goodness of a partition have been proposed, chiefly the Ratio Cut [11] and Normalized Cut [24]. However, minimization of the above measures is an NP-hard problem since it involves searching over all partitions of the nodes. The significance

of spectral clustering partly arises from the fact that it provides a continuous approximation to the above discrete optimization problem [11], [24].

We now describe the RSC algorithm [2]. Denote by  $D = \text{diag}(\hat{d}_1, \dots, \hat{d}_n)$  the diagonal matrix of degrees, where  $\hat{d}_i = \sum_{j=1}^n A_{ij}$ . The normalized (unregularized) symmetric graph Laplacian is defined as

$$L = D^{-1/2}AD^{-1/2}.$$

Regularization is introduced in the following way: Let  $J$  be a constant matrix with all entries equal to  $1/n$ . Then, in regularized spectral clustering one constructs a new adjacency matrix by adding  $\tau J$  to the adjacency matrix  $A$  and computing the corresponding Laplacian. In particular, let

$$A_\tau = A + \tau J,$$

where  $\tau > 0$  is the regularization parameter. The corresponding regularized symmetric Laplacian is defined as

$$L_\tau = D_\tau^{-1/2}A_\tau D_\tau^{-1/2}. \quad (1)$$

Here,  $D_\tau = \text{diag}(\hat{d}_{1,\tau}, \dots, \hat{d}_{n,\tau})$  is the diagonal matrix of ‘degrees’ of the modified adjacency matrix  $A_\tau$ . In other words,  $\hat{d}_{i,\tau} = \hat{d}_i + \tau$ .

The RSC algorithm for finding  $K$  communities is described in Algorithm 1. In order to bring to the forefront the dependence on  $\tau$ , we also denote the RSC algorithm as RSC- $\tau$ . The algorithm first computes  $V_\tau$ , the  $n \times K$  eigenvector matrix corresponding to the  $K$  largest eigenvalues of  $L_\tau$ . The columns of  $V_\tau$  are taken to be orthogonal. The rows of  $V_\tau$ , denoted by  $V_{i,\tau}$ , for  $i = 1, \dots, n$ , corresponds to the nodes in the graph. Clustering the rows of  $V_\tau$ , for example using the  $K$ -means algorithm, provides a clustering of the nodes. We remark that the RSC-0 Algorithm corresponds to the usual spectral clustering algorithm.

---

**Algorithm 1** The RSC- $\tau$  Algorithm [2]

---

**Input :** Laplacian matrix  $L_\tau$ .

**Step 1:** Compute the  $n \times K$  eigenvector matrix  $V_\tau$ .

**Step 2:** Use the  $K$ -means algorithm to cluster the rows of  $V_\tau$  into  $K$  clusters.

---

Our theoretical results assume that the data is randomly generated from a stochastic block model (SBM), which we review in the next subsection. While it is well known that there are real data examples where the SBM fails to provide a good approximation, we believe that the above provides a good playground for understanding the role of regularization in the RSC algorithm. Recent works [2], [10], [23], [6], [14] have used this model, and its variants, to provide a theoretical analyses for various community detection algorithms.

In Chaudhuri et al. [7], the following alternative regularized version of the symmetric Laplacian is proposed:

$$L_{deg,\tau} = D_\tau^{-1/2}AD_\tau^{-1/2}. \quad (2)$$

Here, the subscript *deg* stands for ‘degree’ since the usual Laplacian is modified by adding  $\tau$  to the degree matrix  $D$ . Notice that for the RSC algorithm the matrix  $A$  in the above expression was replaced by  $A_\tau$ .

As mentioned before, we attempt to understand regularization in the framework of the SBM and its extension. We review the SBM in the next section. Using recent results on the concentration of random graph Laplacians [21], we were able to show concentration results in Theorem 4 for the regularized Laplacian in the RSC algorithm. Previous concentration results for the Laplacian (2), as in [7], provide high probability bounds on the spectral norm of the difference of the sample and population regularized Laplacians that depends inversely on  $1/\sqrt{\tau}$ . However, for the regularization (1) we show that the dependence is inverse in  $\tau$ , for large  $\tau$ . We believe that this holds for the regularization (2) as well. We also demonstrate that the eigen gap depends inversely on  $\tau$ , for large  $\tau$ . The benefit of this, along with our improved concentration bounds, is that one can understand regularization by looking at the case where  $\tau$  is large. This results in a very neat criterion for the cluster recovery with the RSC- $\tau$  algorithm.

## 2 The Stochastic Block Model

Given a set of  $n$  nodes, the stochastic block model (SBM), introduced in [12], is one among many random graph models that has communities inherent in its definition. We denote the number of communities in the SBM by  $K$ . Throughout this paper we assume that  $K$  is known. The communities, which represent a partition of the  $n$  nodes, are assumed to be fixed beforehand. Denote these by  $C_1, \dots, C_K$ . Let  $n_k$ , for  $k = 1, \dots, K$ , denote the number of nodes belonging to each of the clusters.

Given the communities, the edges between nodes, say  $i$  and  $j$ , are chosen independently with probability depending the communities  $i$  and  $j$  belong to. In particular, for a node  $i$  belonging to cluster  $C_{k_1}$ , and node  $j$  belonging to cluster  $C_{k_2}$ , the probability of edge between  $i$  and  $j$  is given by

$$P_{ij} = B_{k_1, k_2}.$$

Here, the *block probability matrix*

$$B = ((B_{k_1, k_2})), \quad \text{where } k_1, k_2 = 1, \dots, K$$

is a symmetric full rank matrix, with each entry between  $[0, 1]$ . The  $n \times n$  edge probability matrix  $P = ((P_{ij}))$ , given by (3), represents the population counterpart of the adjacency matrix  $A$ .

Denote  $Z = ((Z_{ik}))$  as the  $n \times K$  binary matrix providing the cluster memberships of each node. In other words, each row of  $Z$  has exactly one 1, with  $Z_{ik} = 1$  if node  $i$  belongs to  $C_k$ . Notice that,

$$P = ZBZ'. \tag{3}$$

Here  $Z'$  denotes the transpose of  $Z$ . Consequently, from (3), it is seen that the rank of  $P$  is also  $K$ .

The population counterpart for the degree matrix  $D$  is denoted by  $\mathcal{D} = \text{diag}(d_1, \dots, d_n)$ , where  $\mathcal{D} = \text{diag}(P\mathbf{1})$ . Here  $\mathbf{1}$  denotes the column vector of all ones. Similarly, the population version of the symmetric Laplacian  $L_\tau$  is denoted by  $\mathcal{L}_\tau$ , where

$$\mathcal{L}_\tau = \mathcal{D}_\tau^{-1/2} P_\tau \mathcal{D}_\tau^{-1/2}.$$

Here  $\mathcal{D}_\tau = \mathcal{D} + \tau I$  and  $P_\tau = P + \tau J$ . The  $n \times n$  matrices  $\mathcal{D}_\tau$  and  $P_\tau$  represent the population counterparts to  $D_\tau$  and  $A_\tau$  respectively. Notice that since  $P$  has rank  $K$ , the same holds for  $\mathcal{L}_\tau$ .

## 2.1 Notation

We use  $\|\cdot\|$  to denote the spectral norm of a matrix. Notice that for vectors this corresponds to the usual  $\ell_2$ -norm. We use  $A'$  to denote the transpose of a matrix, or vector,  $A$ .

For positive  $a_n, b_n$ , we use the notation  $a_n \asymp b_n$  if there exists universal constants  $c_1, c_2 > 0$  so that  $c_1 a_n \leq b_n \leq c_2 a_n$ . Further, we use  $b_n \lesssim a_n$  if  $b_n \leq c_2 a_n$ , for some positive  $c_2$  not depending on  $n$ . The notation  $b_n \gtrsim a_n$  is analogously defined.

The quantities

$$d_{\min, n} = \min_{i=1, \dots, n} d_i, \quad d_{\max, n} = \max_{i=1, \dots, n} d_i$$

denote the minimum and maximum expected degrees of the nodes.

## 2.2 The Population Cluster Centers

We now proceed to define population cluster centers  $\text{cent}_{k, \tau} \in \mathbb{R}^K$ , for  $k = 1, \dots, K$ , for the  $K$  block SBM. These points are defined so that the rows of the eigenvector matrix  $V_{i, \tau}$ , for  $i \in C_k$ , are expected to be scattered around  $\text{cent}_{k, \tau}$ .

Denote by  $\mathcal{V}_\tau$  an  $n \times K$  matrix containing the eigenvectors of the  $K$  largest eigenvalues of the population Laplacian  $\mathcal{L}_\tau$ . As with  $V_\tau$ , the columns of  $\mathcal{V}_\tau$  are also assumed to be orthogonal.

Notice that both  $\mathcal{V}_\tau$  and  $-\mathcal{V}_\tau$  are eigenvector matrices corresponding to  $\mathcal{L}_\tau$ . This ambiguity in the definition of  $\mathcal{V}_\tau$  is further complicated if an eigenvalue of  $\mathcal{L}_\tau$  has multiplicity greater than one. We do away with this ambiguity in the following way: Let  $\mathcal{H}$  denote the set of all  $n \times K$  eigenvector matrices of  $\mathcal{L}_\tau$  corresponding to the top  $K$  eigenvalues. We take,

$$\mathcal{V}_\tau = \arg \min_{H \in \mathcal{H}} \|V_\tau - H\|, \quad (4)$$

where recall that  $\|\cdot\|$  denotes the spectral norm. The matrix  $\mathcal{V}_\tau$ , as defined above, represents the population counterpart of the matrix  $V_\tau$ .

Let  $\mathcal{V}_{i, \tau}$  denote the  $i$ -th row of  $\mathcal{V}_\tau$ . Notice that since the set  $\mathcal{H}$  is closed under the  $\|\cdot\|$  norm, one has that  $\mathcal{V}_\tau$  is also an eigenvector matrix of  $\mathcal{L}_\tau$  corresponding

to the top  $K$  eigenvalues. Consequently, the rows  $\mathcal{V}_{i,\tau}$  are the same across nodes belonging to a particular cluster (See, for example, Rohe et al. [23] for a proof of this fact). In other words, there are  $K$  distinct rows of  $\mathcal{V}_{i,\tau}$ , with each row corresponding to nodes from one of the  $K$  clusters.

Notice that the matrix  $\mathcal{V}_{i,\tau}$  depends on the sample eigenvector matrix  $V_\tau$  through (4), and consequently is a random quantity. However, the following lemma shows that the pairwise distances between the rows of  $\mathcal{V}_{i,\tau}$  are non-random and, more importantly, independent of  $\tau$ .

**Lemma 1.** *Let  $i \in C_k$  and  $i' \in C_{k'}$ . Then,*

$$\|\mathcal{V}_{i,\tau} - \mathcal{V}_{i',\tau}\| = \begin{cases} 0, & \text{if } k = k' \\ \sqrt{\frac{1}{n_k} + \frac{1}{n_{k'}}}, & \text{if } k \neq k' \end{cases}$$

From the above lemma, there are  $K$  distinct rows of  $\mathcal{V}_\tau$  corresponding to the  $K$  clusters. We denote these as  $\text{cent}_{1,\tau}, \dots, \text{cent}_{K,\tau}$ . We also call these the population cluster centers since, intuitively, in an idealized scenario the data points  $V_{i,\tau}$ , with  $i \in C_k$ , should be concentrated around  $\text{cent}_{k,\tau}$ .

### 2.3 Cluster recovery using $K$ -means algorithm

Recall that the RSC- $\tau$  Algorithm 1 works by performing  $K$ -means clustering on the rows of the  $n \times K$  sample eigenvector matrix, denoted by  $V_{i,\tau}$ , for  $i = 1, \dots, n$ . In this section, in particular Corollary 3, we relate the fraction of mis-clustered nodes using the  $K$ -means algorithm to the various parameters in the SBM.

In general, the  $K$ -means algorithm can be described as follows: Assume one wants to find  $K$  clusters, for a given set of data points  $x_i \in \mathbb{R}^K$ , for  $i = 1, \dots, n$ . Then the  $K$ -clusters resulting from applying the  $K$ -means algorithm corresponds to a partition  $\hat{\mathcal{T}} = \{\hat{T}_1, \dots, \hat{T}_K\}$  of  $\{1, \dots, n\}$  that aims to minimize the following objective function over all such partitions:

$$\text{Obj}(\mathcal{T}) = \sum_{k=1}^K \sum_{i \in T_k} \|x_i - \bar{x}_{T_k}\|^2, \quad (5)$$

Here  $\mathcal{T} = \{T_1, \dots, T_K\}$  is a partition  $\{1, \dots, n\}$ , and  $\bar{x}_{T_k}$  corresponds to the vector of component-wise means of the  $x_i$ , for  $i \in T_k$ .

In our situation there is also an underlying true partition of nodes into clusters, given by  $\mathcal{C} = \{C_1, \dots, C_K\}$ . Notice that  $\mathcal{C} = \hat{\mathcal{T}}$  iff there is a permutation  $\pi$  of  $\{1, \dots, K\}$  so that  $C_k = \hat{T}_{\pi(k)}$ , for  $k = 1, \dots, K$ . In general, we use the following measure to quantify the closeness of the outputted partition  $\hat{\mathcal{T}}$  and the true partition  $\mathcal{C}$ : Denote the *clustering error* associated with  $\hat{T}_1, \dots, \hat{T}_K$  as

$$\hat{f} = \min_{\pi} \max_k \frac{|C_k \cap \hat{T}_{\pi(k)}^c| + |C_k^c \cap \hat{T}_{\pi(k)}|}{n_k}. \quad (6)$$

The clustering error measures the maximum proportion of nodes in the symmetric difference of  $C_k$  and  $\hat{T}_{\pi(k)}$ .

In many situations, such as ours, there exists population quantities associated with each cluster around which the  $x_i$ 's are expected to concentrate. Denote these quantities by  $m_1, \dots, m_K$ . In our case,  $m_k = \text{cent}_{k,\tau}$ . If the  $x_i$ 's, for  $i \in C_k$ , concentrate well around  $m_k$ , and the  $m_k$ 's are sufficiently well separated, then it is expected the  $K$ -means algorithm recovers the clusters with small error  $\hat{f}$ .

Denote  $X$  as the  $n \times K$  matrix with  $x_i$ 's as rows. In our case, the  $x_i = V_{i,\tau}$ , and  $X = V_\tau$ . Further, denote as  $M$  the  $n \times K$  matrix with the  $m_k$ 's as rows. In our case,  $M = \mathcal{V}_\tau$ . Recent results on cluster recovery using the  $K$ -means algorithm, as given in Kumar and Kannan [15] and Awasthi and Sheffet [3], provide conditions on  $X$  and  $M$  for the success of  $K$ -means. The following lemma is implied from Theorem 3.1 in Awasthi and Sheffet [3].

**Lemma 2.** *Let  $\delta > 0$  be a small quantity. If for each  $1 \leq k \neq k' \leq K$ , one has*

$$\|m_k - m_{k'}\| \geq \left(\frac{1}{\delta}\right) \sqrt{K} \|X - M\| \left(\frac{1}{\sqrt{n_k}} + \frac{1}{\sqrt{n_{k'}}}\right) \quad (7)$$

*then the clustering error  $\hat{f} = O(\delta^2)$  using the  $K$ -means algorithm.*

**Remark :** In general minimizing the objective function (5) is not computationally feasible. However, the results in [15], [3] can be extended to partitions  $\hat{\mathcal{T}}$  that approximately minimize (5). The condition (7), called the *center separation* condition in [3], provides lower bounds on the pairwise distances between the population cluster centers that depend on the perturbation of data points around the population centers (represented by  $\|X - M\|$ ) and the cluster sizes.

Let

$$1 = \mu_{1,\tau} \geq \dots \geq \mu_{n,\tau}$$

be the eigenvalues of the regularized population Laplacian  $\mathcal{L}_\tau$  arranged in decreasing order. The fact that  $\mu_{1,\tau}$  is 1 follows from standard results on the spectrum of Laplacian matrices (see, for example, [26]). As mentioned in the introduction, in order to control the perturbation of the first  $K$  eigenvectors the eigen gap, given by  $\mu_{K,\tau} - \mu_{K+1,\tau}$ , must be adequately large, as noted in [26], [20], [16]. Since  $\mathcal{L}_\tau$  has rank  $K$  one has  $\mu_{K+1,\tau} = 0$ . Thus the eigen gap is simply  $\mu_{K,\tau}$ . For our  $K$ -block SBM framework the following is an immediate consequence of Lemma 2 and the Davis-Kahan theorem for the perturbation of eigenvectors.

**Corollary 3.** *Let  $\tau \geq 0$  be fixed. For the RSC- $\tau$  algorithm the clustering error, given by (6), is*

$$O\left(\frac{K \|L_\tau - \mathcal{L}_\tau\|^2}{\mu_{K,\tau}^2}\right)$$

*Proof.* Use Lemma 2 with  $m_k = \text{cent}_{k,\tau}$ ,  $X = V_\tau$ ,  $M = \mathcal{V}_\tau$ , and notice that from Lemma 1 that  $\|m_k - m_{k'}\|$  is  $\sqrt{1/n_k + 1/n_{k'}}$ .

Consequently, using  $1/\sqrt{n_k} + 1/\sqrt{n_{k'}} \geq \sqrt{1/n_k + 1/n_{k'}}$  one gets from (7) that if

$$\|V_\tau - \mathcal{V}_\tau\| \leq \frac{\delta}{\sqrt{K}}, \quad (8)$$

for some  $\delta > 0$ , then at most  $O(\delta^2)$  fraction of nodes are misclassified with the RSC- $\tau$  algorithm.

From the Davis-Kahan theorem [5], one has

$$\|V_\tau - \mathcal{V}_\tau\| \lesssim \frac{\|L_\tau - \mathcal{L}_\tau\|}{\mu_{K,\tau}} \quad (9)$$

Consequently, if we take  $\delta = (\sqrt{K}\|L_\tau - \mathcal{L}_\tau\|)/\mu_{K,\tau}$  then relation (8) is satisfied using (9). This proves the corollary.  $\square$

### 3 Improvements through regularization

In this section we will use Corollary 3 to quantify improvements in clustering performance via regularization. If the number of clusters  $K$  is fixed (does not grow with  $n$ ) then the quantity

$$\frac{\|L_\tau - \mathcal{L}_\tau\|}{\mu_{K,\tau}}, \quad (10)$$

in Corollary 3 provides an insight into the role of the regularization parameter  $\tau$ . Clearly, an ideal choice of  $\tau$  would be the one that minimizes (10). Note, however, that this is not practically possible since  $\mathcal{L}_\tau$ ,  $\mu_{K,\tau}$  are not known in advance.

Increasing  $\tau$  will ensure that the Laplacian  $L_\tau$  will be well concentrated around  $\mathcal{L}_\tau$ . This is demonstrated in Theorem 4 below. However, increasing  $\tau$  also has the effect of decreasing the eigen gap, which in this case is  $\mu_{K,\tau}$ , since the population Laplacian becomes more like a constant matrix upon increasing  $\tau$ . Thus the optimum  $\tau$  results from the balancing out of these two competing effects.

Independent of our work, a similar argument for the optimum choice of regularization, using the Davis-Kahan theorem, was given in Qin and Rohe [22] for the regularization proposed in [7]. However, they didn't provide a quantification of the benefit of regularization as given in this section and Section 4.

Theorem 4 provides high-probability bounds on the quantity  $\|L_\tau - \mathcal{L}_\tau\|$  appearing in the numerator of (10). Previous analysis of the regularization (2), in [7], [22], show high-probability bounds on the aforementioned spectral norm that have a  $1/\sqrt{d_{min,n} + \tau}$  dependence on  $\tau$ . However, for large  $\tau$ , the theorem below shows that the behavior is  $\sqrt{d_{max,n}}/(d_{max,n} + \tau)$ . We believe this holds for the regularization (2) as well. Thus, our bounds has a  $1/\tau$  dependence on  $\tau$ , for large  $\tau$ , as opposed to the  $1/\sqrt{\tau}$  dependence shown in [7]. This is crucial since the eigen gap  $\mu_{K,\tau}$  also behaves like  $1/\tau$  for large  $\tau$  which implies that (10) converges to a quantity as  $\tau$  tends to infinity. In Theorem 5 we provide a

bound on this quantity. Our claims regarding improvements via regularization will then follow from comparing this bound with the bound on (10) at  $\tau = 0$ .

**Theorem 4.** *With probability at least  $1 - 2/n$ , for all  $\tau$  satisfying*

$$\max\{\tau, d_{\min,n}\} \geq 32 \log n, \quad (11)$$

*we have*

$$\|L_\tau - \mathcal{L}_\tau\| \leq \epsilon_{\tau,n}. \quad (12)$$

*Here*

$$\epsilon_{\tau,n} = \begin{cases} \frac{10\sqrt{\log n}}{\sqrt{d_{\min,n} + \tau}}, & \text{if } \tau \leq 2d_{\max,n} \\ \frac{10\sqrt{d_{\max,n} \log n}}{d_{\max,n} + \tau/2}, & \text{if } \tau > 2d_{\max,n} \end{cases}$$

We use Theorem 4, along with Corollary 3, to demonstrate improvements from regularization over previous analyses of eigenvector perturbation. Our strategy for this is as follows: Take

$$\delta_{\tau,n} = \frac{\epsilon_{\tau,n}}{\mu_{K,\tau}}$$

Notice that from Corollary 3 and Theorem 4, one gets that with probability at least  $1 - 2/n$ , for all  $\tau$  satisfying (11), the clustering error is  $O(\delta_{\tau,n}^2)$ . Consequently, it is of interest to study the quantity  $\delta_{\tau,n}$  as a function of  $\tau$ . Define,

$$\delta_n = \lim_{\tau \rightarrow \infty} \delta_{\tau,n}. \quad (13)$$

Although we would have ideally liked to study the quantity,

$$\tilde{\delta}_n = \min_{\max\{\tau, d_{\min,n}\} \gtrsim \log n} \delta_{\tau,n}$$

we study  $\delta_n$  since it is easy to characterize as we shall see in Theorem 5 below. Section 5 introduces a data-driven methodology that is based on finding an approximation for  $\tilde{\delta}_n$ .

Before introducing our main theorem quantifying the performance of RSC- $\tau$  for large  $\tau$  we introduce the following definition.

**Definition 1.** *Let  $\{\tau_n, n \geq 1\}$  be a sequence of the regularization parameters. For the  $K$ -block SBM we say that RSC- $\tau_n$  gives consistent cluster estimates if the error (6) goes to 0, with probability tending to 1, as  $n$  goes to infinity.*

Throughout the remainder of the section we consider a  $K$ -block stochastic block model with the following block probability matrix.

$$B = \begin{pmatrix} p_{1,n} & q_n & \cdots & q_n \\ q_n & p_{2,n} & \cdots & q_n \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & q_n & p_{K,n} \end{pmatrix}. \quad (14)$$

The number of communities  $K$  is assumed to be fixed. Without loss, assume that  $p_{1,n} \geq p_{2,n} \dots \geq p_{K,n}$ . We also assume that  $q_n < p_{K,n}$ . Denote  $w_k = n_k/n$ , for  $k = 1, \dots, K$ . The quantity  $w_k$  represents the proportion of nodes belonging to the  $k$ -th community. Throughout this section we assume that  $\{\tau_n : n \geq 1\}$  is a sequence of regularization parameters satisfying,

$$\frac{\left(\sum_{k=1}^K 1/w_k\right) d_{max,n} \log n}{\tau_n} = o(1) \quad (15)$$

Notice that if the cluster sizes are of the same order, that is  $w_k \asymp 1$ , then the above condition simply states that  $\tau_n$  should grow faster than  $d_{max,n} \log n$ .

Denote  $\gamma_{k,n} = n_k(p_{k,n} - q_n)$ . The following is our main result regarding the impact of regularization.

**Theorem 5.** *For the  $K$  block SBM, with block probability matrix (14),*

$$\delta_n \asymp \frac{(\tilde{m}_{1,n} m_{1,n} - m_{2,n})}{m_{1,n}} \sqrt{d_{max,n} \log n}. \quad (16)$$

Here  $\delta_n$  is given by (13) and

$$m_{1,n} = \sum_{k=1}^K \frac{w_k}{\gamma_{k,n}} \quad (17)$$

$$\tilde{m}_{1,n} = \sum_{k=1}^K \frac{1}{\gamma_{k,n}} \quad (18)$$

$$m_{2,n} = \sum_{k=1}^K \frac{w_k}{\gamma_{k,n}^2} \quad (19)$$

Further, let  $\{\tau_n, n \geq 1\}$  satisfy (15). If  $\delta_n$  goes to 0, as  $n$  tends to infinity, then  $RSC\text{-}\tau_n$  gives consistent cluster estimates.

Theorem 5 will be proved in Appendix B. In particular, the following corollary shows that for the stochastic block model regularized spectral clustering would work even when the minimum degree is of constant order. This is an improvement over recent works on unregularized spectral clustering, such as [18], [7], [23], which required the minimum degree to grow at least as fast as  $\log n$ .

**Corollary 6.** *Let the block probability matrix  $B$  be as in (14). Let  $\{\tau_n, n \geq 1\}$  satisfy (15). Then  $RSC\text{-}\tau_n$  gives consistent cluster estimates under the following scenarios:*

i) *For the  $K$ -block SBM if  $w_k \asymp 1$ , for each  $k = 1, \dots, K$ , and*

$$\frac{(p_{K-1,n} - q_n)^2}{p_{1,n}} \text{ grows faster than } \frac{\log n}{n}. \quad (20)$$

ii) For the 2-block SBM if  $p_2 = q$  and

$$\frac{(p_{1,n} - q_n)^2}{w_1 p_{1,n} + w_2 q_n} \text{ grows faster than } \frac{\log n}{n (\min\{w_1, w_2\})^2}. \quad (21)$$

**Remark :** Regime i) deals with the situation that the clusters sizes are of the same order of magnitude. Regime ii), where  $p_{2,n} = q_n$  mimics a scenario where there is only one cluster. This is a generalization of the *planted clique* problem where  $p_{1,n} = 1$  and  $p_{2,n} = q = 1/2$ . For the planted clique problem (21) translates to requiring that  $\min\{w_1, w_2\}$  grow faster than  $\sqrt{\log n}/\sqrt{n}$  for consistent cluster estimates, which is similar to results in [18].

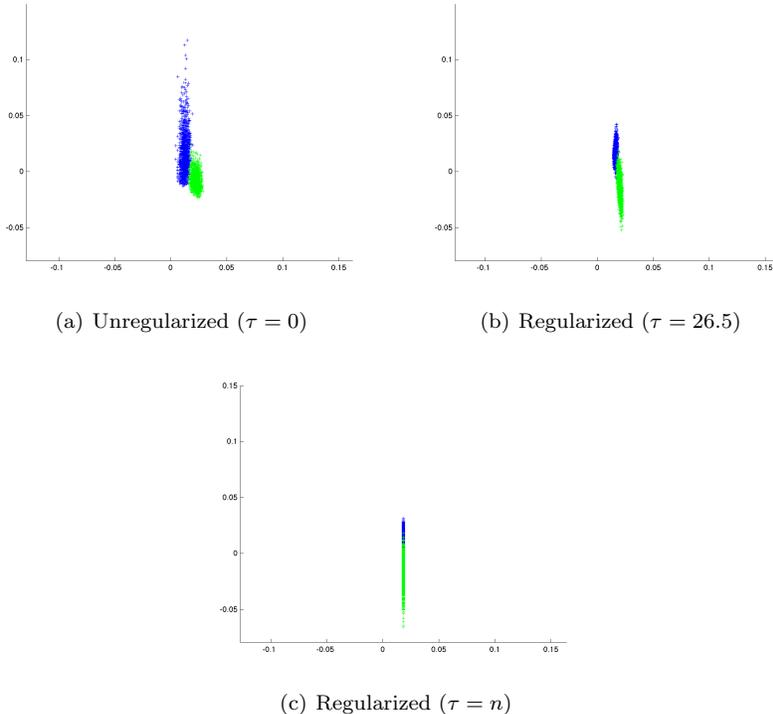


Figure 1: Scatter plot of first two eigenvectors with  $B$  as in (22). The  $x, y$  axes provides values for the first, second eigenvectors respectively. The colors corresponds to the cluster memberships of the nodes. Here the block probability matrix  $B$  is as in (22). Plot a) corresponds to  $\tau = 0$ . b)  $\tau = 26.5$ , selected using our data-driven *DKest* methodology proposed in Section 5. c)  $\tau = n$ .

Notice that in both (20) and (21) the minimum degree could be of constant order. For example, for the two-block SBM if  $q_n, p_{2,n} = O(1/n)$  then the minimum degree is of constant order. In this case ordinary spectral clustering using

the normalized Laplacian would perform poorly. RSC performs better since from (20) it only requires that the larger of the two within block probabilities, that is  $p_{1,n}$ , growing appropriately fast. Figure 1 illustrates this with  $n = 3000$  and edge probability matrix

$$B = \begin{pmatrix} .01 & .0025 \\ .0025 & .003 \end{pmatrix}. \quad (22)$$

The figure provides the scatter plot of the first two eigenvectors of the unregularized and regularized sample Laplacians. Figure a) corresponds to the usual spectral clustering, while plots b) & c) corresponds to RSC- $\tau$ , with  $\tau = 26.5, 3000$  respectively. Here,  $\tau = 26.5$  was selected using our data-driven methodology for selecting  $\tau$  proposed in Section 5. Also,  $\tau = 3000$  was selected as suggested from Theorem 5 and Corollary 6. The fraction of mis-classified are 26%, 4%, 6% for the cases a), b), c) respectively.

From the scatter plots one sees that there is considerably less scattering for the blue points with regularization. This results in improvements in clustering performance. Also, note that the performance in case c), in which  $\tau$  is taken to be very large, is only slightly worse than case b). For case c) there is almost no variation in the first eigenvector, plotted along the  $x$ -axis. This makes sense since the first eigenvector is proportional to  $(\sqrt{\hat{d}_{1,\tau}}, \dots, \sqrt{\hat{d}_{n,\tau}})$  and for large  $\tau$  one has  $\sqrt{\hat{d}_{i,\tau}} \approx \sqrt{\tau}$ .

It may seem surprising that in Corollary 6, claim (20), the smallest within block probability, that is  $p_{K,n}$  does not matter at all. One way of explaining this is that if one can do a good job identifying the top  $K - 1$  highest degree clusters then the cluster with the lowest degree can also be identified simply by eliminating nodes not belonging to this cluster.

## 4 SBM with strong and weak clusters

In many practical situations, not all nodes belong to clusters that can be estimated well. As mentioned in the introduction, these nodes interfere with the clustering of the remaining nodes in the sense that none of the top eigenvectors might discriminate between the nodes that do belong to well-defined clusters. As an example of a real life data set, we consider the political blogs data set, which has two clusters, in Subsection 5.2. With ordinary spectral clustering, the top two eigenvectors do not discriminate between the two clusters (see Figure 2 for explanation). Infact, it is only the third eigenvector that discriminates between the two clusters. This results in bad clustering performance when the first two eigenvectors are considered. However, regularization rectifies this problem by ‘bringing up’ the important eigenvector thereby allowing for much better performance.

We model the above situation – where there are main clusters as well as outlier nodes – in the following way: Consider a stochastic block model, as in

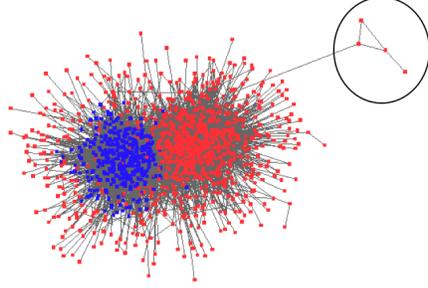


Figure 2: Depiction of the political blog network [1]. Instead of discriminating between the red and blue nodes, the second eigenvector discriminates the small cluster of 4 nodes (circled) from the remaining. This results in bad clustering performance.

(14), with  $K + K_w$  blocks. In particular, let the block probability matrix be given by

$$B = \begin{pmatrix} B_s & B_{sw} \\ B'_{sw} & B_w \end{pmatrix}, \quad (23)$$

where  $B_s$  is a  $K \times K$  matrix with  $(p_{1,n}, \dots, p_{K,n})$  in the diagonal and  $q_n$  in the off-diagonal. Further,  $B_{sw}$ ,  $B_w$  are  $K \times K_w$  and  $K_w \times K_w$  dimensional matrices respectively. In the above  $(K + K_w)$ -block SBM, the top  $K$  blocks corresponds to the well-defined or *strong* clusters, while the bottom  $K_w$  blocks corresponds to less well-defined or *weak* clusters.

We now formalize our notion of strong and weak clusters. The matrix  $B_s$  models the distribution of edges between the nodes belonging to the strong clusters, while the matrix  $B_w$  has the corresponding role for the weak clusters. The matrix  $B_{sw}$  models the interaction between the strong and weak clusters. For ease of analysis, we make the following simplifying assumptions : Assume that  $p_{k,n} = p_n^s$ , for  $k = 1, \dots, K$ , and that the strong clusters  $C_1, \dots, C_K$  have equal sizes, that is, assume  $n_k = n^s$  for  $k = 1, \dots, K$ .

Let  $b_{sw}$  be defined as the maximum of the elements in  $B_{sw}$ , and let  $n^w$  be the number of nodes belonging to a weak cluster. In other words,  $Kn^s + n^w = n$ . We make the following three assumptions:

$$\frac{(p_n^s - q_n)^2}{p_n^s} \text{ grows faster than } \frac{\log n}{n} \quad (24)$$

$$n^w = O(1). \quad (25)$$

$$b_{sw} \lesssim \sqrt{\frac{p_n^s \log n}{n}} \quad (26)$$

Assumption (24) ensures recovery of the strong clusters if there were no nodes belonging to weak clusters (See Corollary 6 or McSherry [18], Corollary

1). Assumption (25) and (26) pertain to the nodes in the weak clusters. In particular, Assumption (25) simply states that the total number of nodes belonging to a weak cluster is constant and does not grow with  $n$ . Assumption (26) states that the density of the edges between the strong and weak clusters, denoted by  $b_{sw}$ , is not too large.

We only assume that the rank of  $B_s$  is  $K$ . Thus, the rank of  $B$  is at least  $K$ . As before, we assume that  $K$  is known and does not grow with  $n$ . The number of weak clusters,  $K_w$ , need not be known and could be as high as  $n^w$ . We do not even place any restriction on the sizes of a weak cluster. Indeed, we even entertain the case that each of the  $K_w$  clusters has one node. Consequently, we are only interested in recovering the strong clusters.

Theorem 7 presents our theorem for the recovery of the  $K$  strong clusters using the RSC- $\tau_n$  Algorithm, with  $\{\tau_n, n \geq 1\}$ , satisfying

$$\frac{np_n^s \log n}{\tau_n} = o(1) \quad (27)$$

In other words, the regularization parameter is taken to grow faster than  $np_n^s \log n$ , where notice that  $np_n^s$  is of the same order of the expected maximum degree of the graph. Let  $\hat{T}_1, \dots, \hat{T}_K$  be the clusters outputted from the RSC- $\tau_n$  Algorithm. Let

$$\hat{f} = \min_{\pi} \max_k \frac{|C_k \cap \hat{T}_{\pi(k)}^c| + |C_k^c \cap \hat{T}_{\pi(k)}|}{n_k},$$

be as in (6). Notice that the clusters  $C_1, \dots, C_K$  do not form a partition of  $\{1, \dots, n\}$ , while the estimates  $\hat{T}_1, \dots, \hat{T}_K$  do. However, since  $n^w$  does not grow with  $n$  this should not make much of a difference.

**Theorem 7.** *Let Assumptions (24), (25) and (26) be satisfied. If  $\{\tau_n, n \geq 1\}$  satisfies (27) then the clustering error  $\hat{f}$  for RSC- $\tau_n$  goes to zero with probability tending to one.*

The theorem is proved in Appendix C. It states that under Assumption (24) – (26) one can get the same results with regularization that one would get if the nodes belonging to the weak clusters weren't present.

Spectral clustering (with  $\tau = 0$ ) may fail under the above assumptions. This is elucidated in Figure 3. Here  $n = 2000$  and there are two strong clusters ( $K = 2$ ) and three weak clusters ( $M = 3$ ). The first 1600 nodes are evenly split between the two strong clusters, with the remaining nodes split evenly between the weak clusters. The matrix  $B_s$  and  $B_w$  are as in (28) and  $B_{sw}$  is a matrix with all entries .015.

$$B_s = \begin{pmatrix} .025 & .015 \\ .015 & .025 \end{pmatrix} \quad B_w = \begin{pmatrix} .007 & .015 & .015 \\ .015 & .0071 & .015 \\ .015 & .015 & .0069 \end{pmatrix}. \quad (28)$$

The nodes in the weak clusters have relatively lower degrees, and consequently, cannot be recovered. Figures 3(a) and 3(b) show the first 3 eigenvectors of the

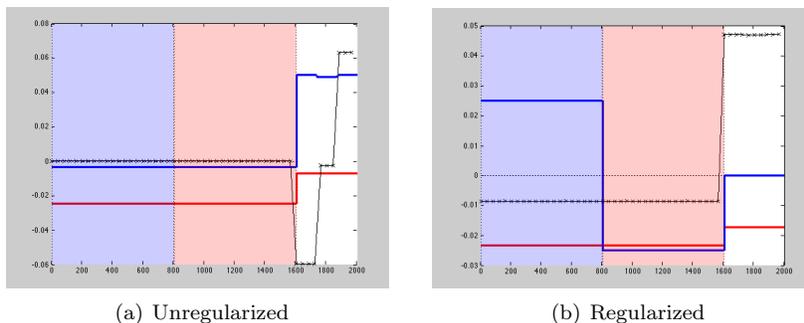


Figure 3: First three population eigenvectors corresponding to  $B_s$  and  $B_w$  in (28). In both plots, the x-axis provides the node indices while the y-axis gives the eigenvector values. The regularization parameter was taken to be  $n$ . The shaded blue and pink regions corresponds to the nodes belonging to the two strong clusters. The solid red line, solid blue line and  $- \times -$  black lines correspond to the first, second and third population eigenvectors respectively.

population Laplacian in the regularized and unregularized cases. We plot the first 3 instead of the first 5 eigenvectors in order to facilitate understanding of the plot. In both cases the first eigenvector is not able to distinguish between the two strong clusters. This makes sense since the first eigenvector of the Laplacian has elements whose magnitude is proportional to square root of the population degrees (see, for example, [26] for a proof of this fact). Consequently, as the population degrees are the same for the two strong clusters, the values for this eigenvector is constant for nodes belonging to the strong clusters.

The situation is different for the second population eigenvector. In the regularized case, the second eigenvector is able to distinguish between these two clusters. However, this is not the case for the unregularized case. From Figure 3(a), not even the third unregularized eigenvector is able to distinguish between the strong and weak clusters. Indeed, it is only the fifth eigenvector that distinguishes between the two strong clusters in the unregularized case.

In Figure 4(a) and 4(b) we show the second sample eigenvector for the two cases in Figure 3(a) and 3(b). Note, we do not show the first sample eigenvector since from Figure 3(a) and 3(b), the corresponding population eigenvectors are not able to distinguish between the two strong clusters. As expected, it is only for the regularized case that one sees that the second eigenvector is able to do a good job in separating the two strong clusters. Running  $K$ -means, with  $k = 2$ , resulted in a mis-classification of 49% of the nodes in the strong clusters in the unregularized case, compared with 16.25% in the regularized case.

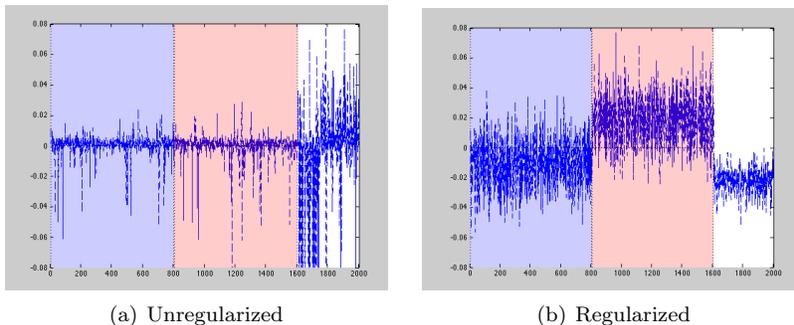


Figure 4: Second sample eigenvector corresponding to situation in Figure 3. As before, in both plots, the x-axis provides the node indices, while the y-axis gives the eigenvector values. As before, the shaded blue and pink regions corresponds to the nodes belonging to the two strong clusters. For plots (a) & (b) the blue line correspond to the second eigenvector of the respective sample Laplacian matrices.

## 5 *DKest* : Data dependent choice of $\tau$

The results Sections 3 and 4 theoretically examined the gains from regularization for large values of regularization parameter  $\tau$ . Those results do not rule out the possibility that intermediate values of  $\tau$  may lead to better clustering performance. In this section we propose a data dependent scheme to select the regularization parameter. We compare it with the scheme in [8] that uses the Girvan-Newman modularity [6]. We use the widely used normalized mutual information criterion (NMI) [2], [27] to quantify the performance of the spectral clustering algorithm in terms of closeness of the estimated clusters to the true clusters.

Our scheme works by directly estimating the quantity in (10) in the following manner: For each  $\tau$  in grid, an estimate  $\hat{\mathcal{L}}_\tau$  of  $\mathcal{L}_\tau$  is obtained using clusters outputted from the RSC- $\tau$  algorithm. In particular, let  $\hat{C}_{1,\tau}, \dots, \hat{C}_{K,\tau}$  be the estimates of the clusters  $C_1, \dots, C_K$  produced from running RSC- $\tau$ . The estimate  $\hat{\mathcal{L}}_\tau$  is taken as the population regularized Laplacian corresponding to an estimated block probability matrix  $\hat{B}$  and clusters  $\hat{C}_{1,\tau}, \dots, \hat{C}_{K,\tau}$ . More specifically, the  $(k_1, k_2)$ -th entry of  $\hat{B}$  is taken as

$$\hat{B}_{k_1, k_2} = \frac{\sum_{i \in \hat{C}_{k_1, \tau}, j \in \hat{C}_{k_2, \tau}} A_{ij}}{|\hat{C}_{k_1, \tau}| |\hat{C}_{k_2, \tau}|} \quad (29)$$

The above is simply the proportion of edges between the nodes in the cluster estimates  $\hat{C}_{k_1, \tau}$  and  $\hat{C}_{k_2, \tau}$ . The following statistic is then considered:

$$DKest_\tau = \frac{\|L_\tau - \hat{\mathcal{L}}_\tau\|}{\mu_K(\hat{\mathcal{L}}_\tau)}, \quad (30)$$

where  $\mu_K(\hat{\mathcal{L}}_\tau)$  denotes the the  $K$ -th smallest eigenvalue of  $\hat{\mathcal{L}}_\tau$ . The  $\tau$  that minimizes the  $DKest_\tau$  criterion is then chosen. Since this criterion provides an estimate of the Davis-Kahan bound, we call it the  $DKest$  criterion.

We compare the above to the scheme that uses Girvan-Newman modularity [6], [19], as suggested in [8]. For a particular  $\tau$  in the grid the Girvan-Newman modularity is computed for the clusters outputted using the RSC- $\tau$  Algorithm. The  $\tau$  that maximizes the modularity value over the grid is then chosen.

Notice that the best possible choice of  $\tau$  would be the one that simply maximizes the NMI over the selected grid. However, this cannot be computed in practice since calculation of the NMI requires knowledge of the true clusters. Nevertheless, this provides a useful benchmark against which one can compare the other two schemes. We call this the ‘oracle’ scheme.

## 5.1 Simulation Results

Figure 5 provides results comparing the three schemes, viz.  $DKest$ , Girvan-Newman and ‘oracle’ schemes. We perform simulations following the pattern of [2]. In particular, for a graph with  $n$  nodes we take the  $K$  clusters to be of equal sizes. The  $K \times K$  block probability matrix is taken to be of the form

$$B = \text{fac} \begin{pmatrix} \beta w_1 & 1 & \dots & 1 \\ 1 & \beta w_2 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ \dots & \dots & 1 & \beta w_K \end{pmatrix}.$$

Here, the vector  $w = (w_1, \dots, w_K)$ , which are the *inside weights*, denotes the relative degrees of nodes within the communities. Further, the quantity  $\beta$ , which is the *out-in ratio*, represents the ratio of the probability of an edge between nodes from different communities to that of probability of edge between nodes in the same community. The scalar parameter  $\text{fac}$  is chosen so that the average expected degree of the graph is equal to  $\lambda$ .

Figure 5 compares the two methods of choosing the best  $\tau$  for various choices of  $n$ ,  $K$ ,  $\beta$ ,  $w$  and  $\lambda$ . In general, we see that the  $DKest$  selection procedure performs at least as well, and in some cases much better, than the procedure that used the Girvan-Newman modularity. The performance of the two methods is much closer when the average degree is small.

## 5.2 Analysis of the Political Blogs dataset

Here we investigate the performance of  $DKest$  on the well studied network of political blogs [1]. The data set aims to study the degree of interaction between liberal and conservative blogs over a period prior to the 2004 U.S Presidential Election. The nodes in the networks are select conservative and liberal blog sites. While the original data set had directed edges corresponding to hyperlinks between the blog sites, we converted it to an undirected graph by connecting two nodes with an edge if there is at least one hyperlink from one node to the other.

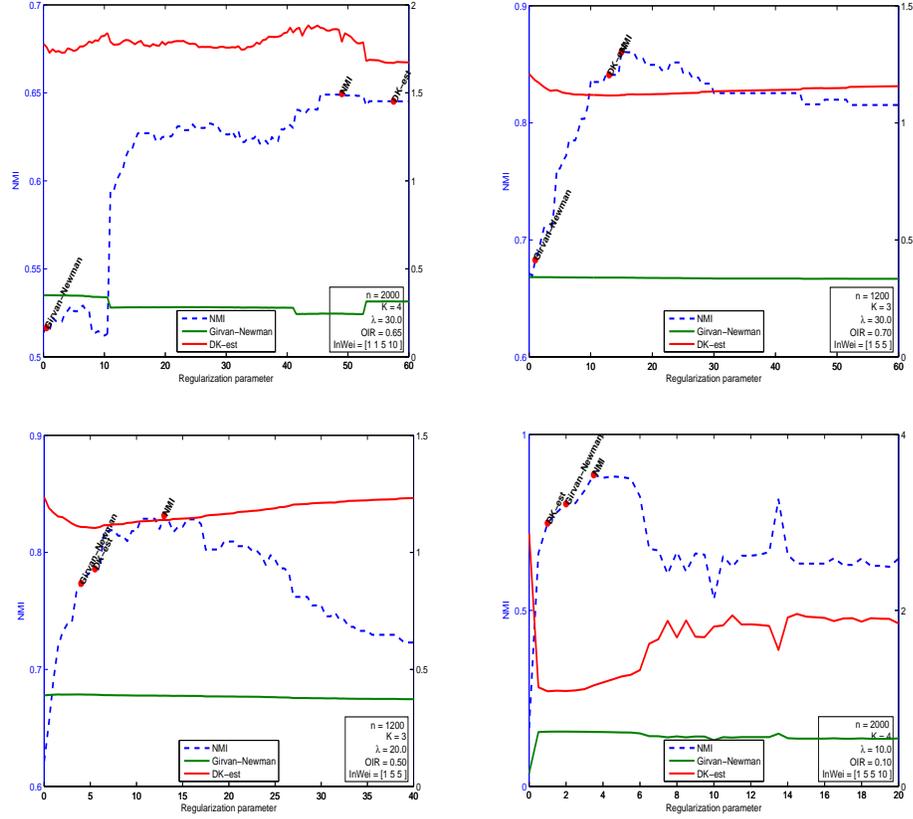


Figure 5: Performance of spectral clustering as a function of  $\tau$  for stochastic block model for  $\lambda$  values of 30, 20 and 10. In the plots we denote  $\beta$  and  $w$  as OIR and  $\text{InWei}$  respectively. The right  $y$ -axis provides values for the Girvan-Newman modularities and  $DKest$  functions, while the left  $y$ -axis provides values for the normalized mutual information (NMI). The 3 labeled dots correspond to values of the NMI at  $\tau$  values which minimizes the  $DKest$ , and maximizes the Girvan-Newman modularity and the NMI. Note, the oracle  $\tau$ , or the  $\tau$  that maximizes the NMI, cannot be calculated in practice.

The data set has 1222 nodes with an average degree of 27. Spectral clustering ( $\tau = 0$ ) resulted in only 51% of the nodes correctly classified as liberal or conservative. The oracle procedure, with  $\tau = 0.5$ , resulted in 95% of the nodes correctly classified. The  $DKest$  procedure selected  $\tau = 2.25$ , with an accuracy of 81%. The Girvan-Newman (GN) procedure, in this case, outperforms the  $DKest$  procedure providing the same accuracy as the oracle procedure. Figure 6 illustrates these findings. As predicted by our theory, the performance becomes

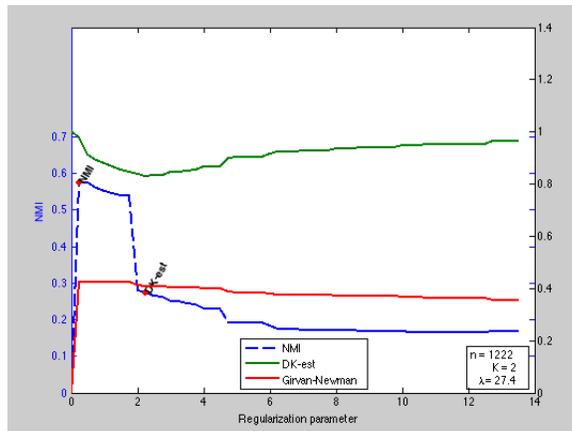


Figure 6: Performance of the three schemes for the political blogs data set [1].

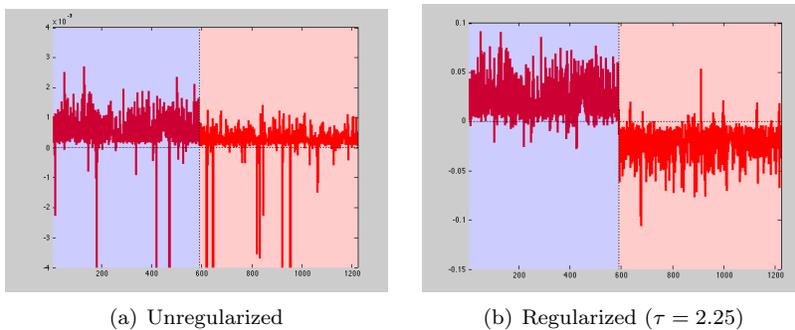


Figure 7: Second eigenvector of the unregularized and regularized Laplacians for the political blogs data set [1]. The shaded blue and pink regions corresponds to the nodes belonging to the liberal and conservative blogs respectively.

insensitive for large  $\tau$ . In this case 70% of the nodes are correctly clustered for large  $\tau$ .

We remark that the *DKest* procedure does not perform as well as the GN procedure most likely because our estimate  $\hat{\mathcal{L}}_\tau$  in (30) assumes that the data is generated from an SBM, which is a poor model for the data due to the large heterogeneity in the node degrees. A better model for the data would be the degree corrected stochastic block model (D-SBM) proposed by Karrer and Newman [14]. If we use D-SBM based estimates in *DKest* then the selection of  $\tau$  matches that of the GN Newman and the oracle procedure. See Section 6 for a discussion on this.

The results of Section 4 also explain why unregularized spectral clustering performs badly (see Figure 2). The first eigenvector in both cases (regularized

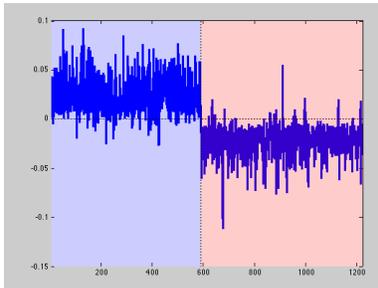


Figure 8: Third eigenvector of the unregularized Laplacian.

and unregularized) does not discriminate between the two clusters. In Figure 7, we plot the second eigenvector of the regularized and unregularized Laplacians. The second eigenvector is able to discriminate between the clusters in the regularized case, while it fails to do so in without regularization. Indeed, it is only the third eigenvector in the unregularized case that distinguishes between the clusters, as shown in Figure 8.

## 6 Discussion

The paper provides a theoretical justification for regularization. In particular, we show why choosing a large regularization parameter can lead to good results. The paper also partly explains empirical findings in Amini et al. [2] showing that the performance of regularized spectral clustering becomes insensitive for larger values of regularization parameters. It is unclear at this stage whether the benefits of regularization, resulting from the trade-offs between the eigen gap and the concentration bound, hold for the regularization in [7], [22] as they hold for the regularization in Amini et al. [2] (as demonstrated in Sections 3 and 4).

Even though our theoretical results focus on larger values of the regularization parameter it is very likely that intermediate values of  $\tau$  produce better clustering performance. Consequently, we propose a data-driven methodology for choosing the regularization parameter. We hope to quantify theoretically the gains from using intermediate values of the regularization parameter in a future work.

For the extension of the SBM proposed in Section 4, if the rank of  $B$ , given by (23), is  $K$  then the model encompasses specific degree-corrected stochastic block models (D-SBM) [14] where the edge probability matrix takes the form

$$P = \Theta Z B Z' \Theta.$$

Here  $\Theta = \text{diag}(\theta_1, \dots, \theta_n)$  models the heterogeneity in the degrees. In particular, consider a  $K$ -block D-SBM with  $0 < \theta_i \leq 1$ , for each  $i$ . Assume that  $\theta_i = 1$  for the most of the nodes. Take the nodes in the strong clusters to be those with  $\theta_i = 1$ . The nodes in the strong clusters are associated to one of  $K$  clusters

depending on the cluster they belong to in the D-SBM. The remaining nodes are taken to be in the weak clusters. Assumptions (25) and (26) puts constraints on the  $\theta_i$ 's which allows one to distinguish between the strong clusters via regularization. It would be interesting to investigate the effect of regularization in more general versions of the D-SBM, especially where there are high as well as low degree nodes.

The *DKest* methodology for choosing the regularization parameter works by providing estimates of the population Laplacian assuming that the data is drawn from an SBM. From our simulations, it is seen that the performance of *DKest* does not change much if we take the matrix norm in the numerator of (30) to be the Frobenius norm, which is much faster to compute.

It is seen that the performance of *DKest* improves for the political blogs data set by taking  $\hat{\mathcal{L}}_\tau$  to be the estimate assuming that the data is drawn from the more flexible D-SBM. Indeed, if we take  $\hat{\mathcal{L}}_\tau$  to be such an estimate then the performance of *DKest* is seen to be as good as the oracle scheme (and the GN scheme) for this data set. We describe how we construct this estimate in Appendix D.

## Acknowledgments

This paper is supported in part by NSF grants DMS-1228246 and DMS-1160319 (FRG), ARO grant W911NF-11-1-0114, NHGRI grant 1U01HG007031-01 (ENCODE), and the Center of Science of Information (CSoI), a US NSF Science and Technology Center, under grant agreement CCF-0939370. A. Joseph would like to thank Sivaraman Balakrishnan and Puramrita Sarkar for some very helpful discussions, and also Arash A. Amini for sharing the code used in the work [2].

## A Analysis of SBM with $K$ blocks

Throughout this section we assume that we have samples from a  $K$  block SBM. Denote the sample and population regularized Laplacian as  $L_\tau, \mathcal{L}_\tau$  respectively. For ease of notation, we remove the subscript  $\tau$  from the various matrices such as  $L_\tau, \mathcal{L}_\tau, A_\tau, D_\tau, \mathcal{D}_\tau$ . We also remove the subscript  $\tau$  in the  $\hat{d}_{i,\tau}, d_{i,\tau}$ 's and denote these as  $\hat{d}_i, d_i$  respectively. However, in some situations we may need to refer to these quantities at  $\tau = 0$ . In such cases, we make this clear by writing them as  $\hat{d}_{i,0}$ , for  $i = 1, \dots, n$  and  $d_{i,0}$  for  $i = 1, \dots, n$ .

We need probabilistic bounds on the weighed sum of Bernoulli random variables. The following lemma is proved in [13].

**Lemma 8.** *Let  $W_j, 1 \leq j \leq N$  be  $N$  independent Bernoulli( $r_j$ ) random variables. Furthermore, let  $\alpha_j, 1 \leq j \leq N$  be non-negative weights that sum to 1 and let  $N_\alpha = 1/\max_j \alpha_j$ . Then the weighted sum  $\hat{r} = \sum_j \alpha_j W_j$ , which has mean given by  $r^* = \sum_j \alpha_j r_j$ , satisfies the following large deviation inequalities. For any  $r$  with  $0 < r < r^*$ ,*

$$P(\hat{r} < r) \leq \exp\{-N_\alpha D(r||r^*)\} \quad (31)$$

and for any  $\tilde{r}$  with  $r^* < \tilde{r} < 1$ ,

$$P(\hat{r} > \tilde{r}) \leq \exp\{-N_\alpha D(\tilde{r}||r^*)\} \quad (32)$$

where  $D(r||r^*)$  denotes the relative entropy between Bernoulli random variables of success parameters  $r$  and  $r^*$ .

The following is an immediate corollary of the above.

**Corollary 9.** *Let  $W_j$  be as in Lemma 8. Let  $\beta_j$ , for  $j = 1, \dots, N$  be non-negative weights, and let*

$$W = \sum_{j=1}^N \beta_j W_j.$$

Then,

$$P(W - E(W) > \delta) \leq \exp\left\{-\frac{1}{2 \max_j \beta_j} \frac{\delta^2}{(E(W) + \delta)}\right\} \quad (33)$$

and

$$P(W - E(W) < -\delta) \leq \exp\left\{-\frac{1}{2 \max_j \beta_j} \frac{\delta^2}{E(W)}\right\} \quad (34)$$

*Proof.* Here we use the fact that

$$D(r||r^*) \geq (r - r^*)^2/(2r), \quad (35)$$

for any  $0 < r, r^* < 1$ . We prove (33). The proof of (34) is similar. The event under consideration may be written as

$$\{\hat{r} - r^* > \tilde{\delta}\},$$

where  $\hat{r} = W/\sum_j \beta_j$ ,  $r^* = E(W)/\sum_j \beta_j$  and  $\tilde{\delta} = \delta/\sum_j \beta_j$ . Correspondingly, using Lemma 8 and (35), one gets that

$$P(W - E(W) > \delta) \leq \exp\left\{-\frac{\sum_j \beta_j}{\max_j \beta_j} \frac{\tilde{\delta}^2}{2(r^* + \tilde{\delta})}\right\}.$$

Substituting the values of  $\tilde{\delta}$  and  $r^*$  results in bound (33).  $\square$

The following lemma provides high probability bounds on the degree. Let  $\tau_{min} = \max\{d_{min,n}, c \log n\}$  and  $\delta_{i,c} = \max\{d_{i,0}, c \log n\}$ .

**Lemma 10.** *On a set  $E_1$  of probability at most  $1 - 2/n^{c_1-1}$ , one has*

$$|\hat{d}_{i,\tau} - d_{i,\tau}| \leq c_2 \sqrt{\delta_{i,c} \log n} \quad \text{for each } i = 1, \dots, n.,$$

where  $c_1 = .5c_2^2/(1 + c_2/\sqrt{c})$ .

*Proof.* Use the fact that  $\hat{d}_{i,\tau} - d_{i,\tau} = \hat{d}_{i,0} - d_{i,0}$ , and

$$P(|\hat{d}_{i,0} - d_{i,0}| \leq c_2 \sqrt{\delta_{i,c} \log n} \quad \forall i) \leq \sum_{i=1}^n P(|\hat{d}_{i,0} - d_{i,0}| \leq c_2 \sqrt{\delta_{i,c} \log n})$$

Notice that  $\hat{d}_{i,0} = \sum_{j=1}^n A_{ij}$ . Apply Corollary 9 with  $\beta_j = 1$  and  $W_j = A_{ij}$ , and  $\delta = c_2 \sqrt{\tau_{min} \log n}$  to bound each term in the sum of the right side of the above equation.

The error exponent can be bounded by,

$$2n \exp \left\{ -\frac{1}{2} \frac{\delta^2}{E(W) + \delta} \right\}. \quad (36)$$

We claim that,

$$E(W) + \delta \leq (1 + c_2/\sqrt{c})\delta_{i,c}. \quad (37)$$

Substituting the above bound in the error exponent (36) will complete the proof.

To see the claim, notice that  $E(W) = d_{i,0}$ . Now, consider the case  $d_{i,0} \geq c \log n$ . In this case,  $\delta_{i,c} = d_{i,0}$  and  $\log n < d_{i,0}/c$ . Correspondingly,  $E(W) + \delta$  is at most  $d_{i,0}(1 + c_2/\sqrt{c})$ .

Next, consider the case  $d_{i,0} < c \log n$ . In this case  $\delta_{i,c} = \tau_{min}$ , which is  $c \log n$ . Consequently,

$$E(W) + \delta \leq c \log n + c_2 \sqrt{c} \log n.$$

The right side of the above can be bounded by  $(1 + c_2/\sqrt{c})(c \log n)$ . This proves the claim.  $\square$

## A.1 Concentration of Laplacian

Below we provide the proof of Theorem 4. Throughout this section we assume that the quantities  $c, c_2$  appearing in Lemma 10 are given by  $c = 32$  and  $c_2 = 2\sqrt{2}$ . Notice that this makes  $c_1 > 2$ , where  $c_1$  as in Lemma 10.

From Lemma 10, with probability at least  $1 - n^{-1}$ ,

$$\max_i |\hat{d}_i - d_i|/d_i \leq \max_i c_2 \sqrt{\delta_{i,c} \log n}/d_i$$

We claim that the right side of the above is at most  $1/2$ . To see this notice that

$$\begin{aligned} \sqrt{\delta_{i,c} \log n}/d_i &\leq \sqrt{\delta_{i,c} \log n}/\delta_{i,c} \\ &= \sqrt{\log n}/\sqrt{\delta_{i,c}} \\ &\leq 1/\sqrt{c} \end{aligned}$$

Here the first inequality follows from noting that  $d_i = d_{i,0} + \tau$ , which is at most  $\max\{d_{i,0}, c \log n\}$ , using  $\tau \geq c \log n$ . The third inequality follows from using  $\delta_{i,c} \geq c \log n$ . Consequently,  $\max_i |\hat{d}_i - d_i|/d_i \leq 1/2$  using  $c_2 = 2\sqrt{2}$  and  $c = 32$ .

*Proof of Theorem 4.* Our proof has parallels with the proof in [21]. Write  $\tilde{L} = \mathcal{D}^{-1/2} A \mathcal{D}^{-1/2}$ . Then,

$$\|L - \mathcal{L}\| \leq \|L - \tilde{L}\| + \|\tilde{L} - \mathcal{L}\|.$$

We first bound  $\|L - \tilde{L}\|$ . Let  $F = D^{1/2} \mathcal{D}^{-1/2}$ . Then  $\tilde{L} = FLF$ . Correspondingly,

$$\begin{aligned} \|L - \tilde{L}\| &\leq \|L - FL\| + \|FL - \tilde{L}\| \\ &\leq \|I - F\| \|L\| + \|F\| \|L\| \|I - F\| \\ &\leq \|I - F\| (2 + \|I - F\|) \end{aligned} \quad (38)$$

Notice that

$$F - I = (I + (D - \mathcal{D}) \mathcal{D}^{-1})^{1/2} - I.$$

Further, using  $\max_i |\hat{d}_i - d_i|/d_i \leq 1/2$ , and the fact that  $\sqrt{1+x} - 1 \leq x$  for  $x \in [-3/4, 3/4]$ , as in [21], one gets that

$$\|F - I\| \leq c_2 \frac{\max_i \sqrt{\delta_{i,c} \log n}}{d_i}$$

with high probability. Consequently, using (38), one gets that

$$\|L - \tilde{L}\| \leq c_2 \max_i \frac{\sqrt{\delta_{i,c} \log n}}{d_i} \left( 2 + c_2 \max_i \frac{\sqrt{\delta_{i,c} \log n}}{d_i} \right) \quad (39)$$

with probability at least  $1 - 1/n^{c_1-1}$ .

$$\max_i \frac{\sqrt{\delta_{i,c}}}{d_i} \leq \tilde{\epsilon}_{\tau,n} = \begin{cases} \frac{1}{\sqrt{d_{\min,n} + \tau}}, & \text{if } \tau \leq 2d_{\max,n} \\ \frac{\sqrt{d_{\max,n}}}{d_{\max,n} + \tau/2}, & \text{if } \tau > 2d_{\max,n} \end{cases}$$

To see this notice, that  $\delta_{i,c} \leq d_{i,0} + \tau = d_i$ , using  $\max\{\tau, d_{i,0}\} \geq c \log n$ . Consequently,  $\sqrt{\delta_{i,c}}/d_i \leq 1/\sqrt{d_{i,0} + \tau}$ , which is at most  $1/\sqrt{d_{\min,n} + \tau}$ .

Further,

$$\max_i \frac{\sqrt{\delta_{i,c}}}{d_i} \leq \frac{\sqrt{d_{\max,n}}}{d_{\max,n} + \tau}$$

for  $\tau > d_{\max,n}$ . This is atmost  $\tilde{\epsilon}_{\tau,n}$  for  $\tau > d_{\max,n}$ .

Consequently, from (39), one gets that

$$\|L - \tilde{L}\| \leq c_2 \tilde{\epsilon}_{\tau,n} \sqrt{\log n} (2 + c_2/\sqrt{c}) \quad (40)$$

with probability at least  $1 - 1/n^{c_1-1}$ .

Next, we bound  $\|\tilde{L} - \mathcal{L}\|$ . We get high probability bounds on this quantity using results in [21], [17]. In particular, as in [21],

$$\tilde{L} - \mathcal{L} = \sum_{i \leq j} Y_{ij},$$

where  $Y_{ij} = \mathcal{D}^{-1/2} X_{ij} \mathcal{D}^{-1/2}$ , with

$$X_{ij} = \begin{cases} (A_{ij} - P_{ij}) (e_i e_j^T + e_j e_i^T), & \text{if } i \neq j \\ (A_{ij} - P_{ij}) e_i e_i^T & \text{if } i = j \end{cases}.$$

Further,  $\|Y_{ij}\| \leq 1/(d_{\min,n} + \tau)$ . Let  $\sigma^2 = \|\sum_{i \leq j} E(Y_{ij}^2)\|$ . We claim that  $\sigma^2 \leq \tilde{\epsilon}_{\tau,n}^2$ . As in [21], page 15, notice that,

$$\sum_{i \leq j} E(Y_{ij}^2) = \sum_{i=1}^n \frac{1}{d_{i,0} + \tau} \left( \sum_{j=1}^n \frac{P_{ij}(1 - P_{ij})}{d_{j,0} + \tau} \right) e_i e_i^T. \quad (41)$$

Clearly,

$$\left( \sum_{j=1}^n \frac{P_{ij}(1 - P_{ij})}{d_{j,0} + \tau} \right) \leq \frac{d_{i,0}}{d_{\min,n} + \tau}.$$

Consequently, for each  $i$  the right side of (41) is at most  $1/(d_{\min,n} + \tau)$  leading to the fact that  $\sigma^2 \leq 1/(d_{\min,n} + \tau)$ .

For  $\tau > 2d_{\max,n}$  we can get improvements in the bound for  $\sigma^2$ . By using the fact that  $d_{j,0} + \tau > d_{\max,n} + \tau/2$  for  $\tau > 2d_{\max,n}$ , one gets that

$$\left( \sum_{j=1}^n \frac{P_{ij}(1 - P_{ij})}{d_{j,0} + \tau} \right) \leq \frac{d_{i,0}}{d_{\max,n} + \tau/2}.$$

for  $\tau > 2d_{\max,n}$ . Consequently, using  $d_{i,0}/(d_{i,0} + \tau) \leq d_{\max,n}/(d_{\max,n} + \tau)$ , one gets that  $\sigma^2 \leq d_{\max,n}/(d_{\max,n} + \tau/2)^2$  for  $\tau > 2d_{\max,n}$ .

Applying Corollary 4.2 in [17] one gets

$$P\left(\|\tilde{L}_0 - \mathcal{L}_0\| \geq t\right) \leq n e^{-t^2/2\sigma^2}.$$

Consequently, with probability at least  $1 - 1/n^{c_1-1}$  one has,

$$\|\tilde{L} - \mathcal{L}\| \leq \sqrt{\frac{2c_1 \log n}{d_{\min,n}}}.$$

Thus, with probability at least  $1 - 1/n^{c_1-1}$ , one has

$$\|\tilde{L} - \mathcal{L}\| \leq \sqrt{2c_1 \log n} \tilde{\epsilon}_{\tau,n}. \quad (42)$$

As a result, combining (40) and (42), one gets that with probability at least  $1 - 2/n^{c_1-1}$ , one has

$$\|L_\tau - \mathcal{L}_\tau\| \leq \sqrt{\log n} \tilde{\epsilon}_{\tau,n} [\sqrt{2c_1} + c_2 (2 + (c_2/\sqrt{c}))]$$

Substituting the values of  $c_2$ ,  $c$ , and noting that  $c_1 > 2$  one gets the expression in the theorem.  $\square$

## A.2 Proof of Lemma 1

Notice that the population regularized Laplacian  $\mathcal{L}_\tau$  corresponds to the population Laplacian of an ordinary stochastic block model with block probability matrix

$$B_\tau = B + vv',$$

where  $v = (\sqrt{\tau/n})\mathbf{1}$ . Correspondingly, we can use the following facts of the population eigenvectors and eigenvalues given for a SBM.

Let  $Z$  be the community membership matrix, that is, the  $n \times K$  matrix with entry  $(i, k)$  being 1 if node  $i$  belongs to cluster  $C_k$ . The following is proved in [23]:

1. Let  $R = \mathcal{D}_\tau^{-1}$ . Then, the non-zero eigen values of  $\mathcal{L}_\tau$  are the same as that of

$$B_{\text{eig}} = B_\tau(Z'RZ), \quad (43)$$

or equivalently,  $\tilde{B}_{\text{eig}} = (Z'RZ)^{1/2}B_\tau(Z'RZ)^{1/2}$ .

2. Define  $\mu = R^{1/2}Z(Z'RZ)^{-1/2}$ . Let,

$$\tilde{B}_{\text{eig}} = H\Lambda H^T,$$

where the right side of the above gives the singular value decomposition of the matrix on the right. Then the eigenvectors of  $\mathcal{L}_\tau$  are given by  $\mu H$ .

Further, since in the stochastic block model the expected node degrees are the same for all nodes in a particular cluster, one can write  $R^{1/2}Z = ZQ$ , where  $Q^{-2}$  is the  $K \times K$  diagonal matrix of population degrees of nodes in a particular community. Consequently, one sees that

$$\mu H = Z(Z^T Z)^{-1/2}H.$$

Lemma 1 follows from noting that

$$\mu H(\mu H)^T = Z(Z^T Z)^{-1}Z^T$$

and the fact that  $(Z^T Z)^{-1} = \text{diag}(1/n_1, \dots, 1/n_k)$ .

## B Proof of Theorem 5

We first prove (16). Recall that  $\delta_{\tau,n}$  is the limit of  $\epsilon_{\tau,n}/\mu_{K,\tau}$ , as  $\tau \rightarrow \infty$ . Now  $\tau\epsilon_{\tau,n}$  converges to  $20\sqrt{d_{\max,n}} \log n$ . Consequently, we now show that

$$\lim_{\tau \rightarrow \infty} \frac{1}{\tau\mu_{K,\tau}} \asymp \frac{\tilde{m}_{1,n}m_{1,n} - m_{2,n}}{m_{1,n}}. \quad (44)$$

Recall that  $\mu_{K,\tau}$  is the  $K$ -th smallest eigenvalue of  $B_{\text{eig}}$  (43). Now,

$$\mu_{K,\tau} \asymp \frac{1}{\text{trace}(B_{\text{eig}}^{-1})}.$$

The above follows from noting that  $\mu_{K,\tau}$  is also equal to the inverse of the largest eigenvalue of  $B_{eig}^{-1}$ , and the fact that the latter is  $\asymp \text{trace}(B_{eig}^{-1})$ , as  $K$  is fixed. We now proceed to show that  $\text{trace}(B_{eig}^{-1})/\tau$  converges to a quantity that is of the same order of magnitude as the right side of (44). This will prove (16).

Recall that the block probability matrix  $B$  is given by (14). We first consider the case that  $q_n = 0$ , that is, there is no interaction between the clusters. Notice,

$$B_{eig}^{-1} = F^{-1}(B + vv')^{-1},$$

where

$$F^{-1} = \text{diag}\left(\frac{\gamma_1 + \tau}{n_1}, \dots, \frac{\gamma_K + \tau}{n_k}\right).$$

Here, for convenience, we remove the subscript  $n$  from quantities such as  $\gamma_{i,n}$ . Using Sherman-Morrison formula

$$(B + vv')^{-1} = B^{-1} - \frac{(B^{-1}v)(B^{-1}v)'}{1 + v'B^{-1}v}$$

One sees that,  $B^{-1}v = \sqrt{\tau/n}(1/p_1, \dots, 1/p_K)'$ . Correspondingly,

$$v'B^{-1}v = \frac{\tau}{n} \sum_i 1/p_i = \tau m_{1,n},$$

using  $q_n = 0$ . Further, the diagonal entries of the matrix  $(B^{-1}v)(B^{-1}v)'$  can be written as

$$\frac{\tau}{n} \text{diag}(1/p_1^2, \dots, 1/p_K^2).$$

We need the trace of  $B_{eig}^{-1}$ . Using the above, one sees that

$$\text{trace}(B_{eig}^{-1}) = \sum_k \frac{\gamma_k + \tau}{\gamma_k} - \frac{\tau m_{1,n} + \tau^2 m_{2,n}}{1 + \tau m_{1,n}}.$$

Since  $K$  is fixed, we have,

$$\text{trace}(B_{eig}^{-1}) \asymp \tau \tilde{m}_{1,n} - \frac{\tau m_{1,n} + \tau^2 m_{2,n}}{1 + \tau m_{1,n}}. \quad (45)$$

Thus, as  $\tau \rightarrow \infty$ , one gets that,

$$\frac{\text{trace}(B_{eig}^{-1})}{\tau} \text{ converges to } \tilde{m}_{1,n} - m_{2,n}/m_{1,n}.$$

The right side of the above is positive, as  $\tilde{m}_{1,n} m_{1,n} \geq m_{2,n}$ , for  $K > 1$ .

Now consider the  $K$  block model with off-diagonal elements of  $B$  equal to  $q$ . Notice that

$$B_\tau = B_0 + \tilde{v}(\tilde{v})^T,$$

where  $B_0 = \text{diag}(p_1 - q, \dots, p_K - q)$  and  $\tilde{v} = \sqrt{\tilde{\tau}/n}\mathbf{1}$ , where  $\tilde{\tau} = \tau + nq$ . Thus applying the above result for the diagonal block model one gets that if  $\tau$  tends

to infinity, the quantity  $\text{trace}(B_{\text{eig}}^{-1})/\tau$  converges to  $\tilde{m}_{1,n} - m_{2,n}/m_{1,n}$ , where here  $\gamma_k = n_k(p_k - q)$ . This proves (16).

We now prove that RSC- $\tau_n$  provides consistent cluster estimates for  $\{\tau_n, n \geq 1\}$  satisfying (15). We need to show that  $\epsilon_{\tau_n, n}/\mu_{K, \tau_n}$  goes to zero.

First, notice that  $\tau_n \epsilon_{\tau_n, n} \lesssim \sqrt{d_{\max, n} \log n}$ . Consequently, from the above, we need to show that  $\text{trace}(B_{\text{eig}}^{-1})\sqrt{d_{\max, n} \log n}/\tau_n$  is  $o(1)$  if  $\delta_n = o(1)$ . From (45) one has

$$\frac{\text{trace}(B_{\text{eig}}^{-1})\sqrt{d_{\max, n} \log n}}{\tau_n} \asymp \sqrt{d_{\max, n} \log n} \left[ \frac{\tilde{m}_{1,n} - m_{1,n}}{1 + \tau_n m_{1,n}} + (\tau_n m_{1,n}) \frac{\tilde{m}_{1,n} - m_{2,n}/m_{1,n}}{1 + \tau_n m_{1,n}} \right]$$

The second term is bounded by  $\delta_n$ , which, by assumption, goes to zero. The first term is bounded by  $\sqrt{d_{\max, n} \log n} \tilde{m}_{1,n}/(m_{1,n} \tau_n)$ . Noting that  $\tilde{m}_{1,n}/m_{1,n} \lesssim \sum_k 1/w_k$ , one gets that the second terms also goes to 0, as  $\tau_n$  satisfies (15).

## B.1 Proof of Corollary 6

For the  $K$ -block SBM, let  $r_K = \gamma_{K,n}/\gamma_{K-1,n}$ . Notice that  $r_K \asymp (p_{K-1} - q)/(p_K - q)$  using  $w_k \asymp 1$ . Use the fact that  $m_{1,n} = (1/\gamma_{K,n})(w_K + O(r_K))$ ,  $\tilde{m}_{1,n} = (1/\gamma_{K,n})(1 + O(r_K))$  and  $m_{2,n} = (1/\gamma_{K,n}^2)(w_K + O(r_K))$ , to get that

$$\frac{(\tilde{m}_{1,n} m_{1,n} - m_{2,n})}{m_{1,n}} = O(1/\gamma_{K-1,n}).$$

Consequently,  $\delta_n = O(\sqrt{d_{\max, n} \log n}/\gamma_{K-1,n})$ . The proof of claim (20) is completed by noting that  $\gamma_{K-1,n} \asymp n(p_{K-1} - q)$  and  $d_{\max, n} \asymp n p_{1,n}$ .

For the 2-block SBM we show that

$$\delta_n \asymp \frac{\sqrt{d_{\max, n} \log n}}{w_1 w_2 [(p_{1,n} + p_{2,n})/2 - q_n]}. \quad (46)$$

Expression (46) follows from using (16) and noting that

$$\frac{(\tilde{m}_{1,n} m_{1,n} - m_{2,n})}{m_{1,n}} = \frac{1}{w_2 \gamma_{1,n} + w_1 \gamma_{2,n}}$$

for the two-block model. It is seen that

$$w_2 \gamma_{1,n} + w_1 \gamma_{2,n} = 2n w_1 w_2 [(p_{1,n} + p_{2,n})/2 - q_n].$$

Notice that  $w_1 w_2 \asymp \min\{w_1, w_2\}$ . Consequently, (21) follows from noting that when  $p_{2,n} = q_n$  then  $d_{\max, n} = n(w_1 p_{1,n} + w_2 q_n)$ .

## C Proof of Results in Section 4

In this section we provide the proof Theorem 7, along with Lemmas 11 and 12 required in proving the theorem.

## C.1 Proof of Theorem 7

Denote  $C^w$  as the set of nodes belonging to the weak clusters. We club all the nodes belonging to the weak clusters into the cluster  $C_K$  and call this combined cluster as  $\tilde{C}_K$ , that is  $\tilde{C}_K = C_K \cup C^w$ . For consistency of notation, let  $\tilde{C}_k = C_k$ , for  $1 \leq k \leq K-1$ , and let  $\tilde{n}_k = |\tilde{C}_k|$ , for  $k = 1, \dots, K$ .

Denote

$$\tilde{f} = \min_{\pi} \max_k \frac{|\tilde{C}_k \cap \hat{T}_{\pi(k)}^c| + |\tilde{C}_k^c \cap \hat{T}_{\pi(k)}|}{\tilde{n}_k}.$$

It is not hard to see that,

$$\hat{f} \leq \left(1 + \frac{n^w}{n^s}\right) \tilde{f} + \frac{n^w}{n^s}.$$

Consequently, a demonstration the  $\tilde{f}$  goes to zero, along with the fact that  $n^w = O(1)$ , will show that  $\hat{f}$  goes to zero.

We now show that  $\tilde{f}$  goes to zero with high probability. For a given assignment of nodes in one of the  $K + K_w$  clusters we denote  $L_{\tau}, \mathcal{L}_{\tau}$  to be the sample, population regularized Laplacians respectively. Further, let  $\tilde{\mathcal{L}}_{\tau}$  be the population regularized Laplacian of a  $K + 1$ -block SBM constructed from clusters  $C_1, \dots, C_K$  and  $C^w$ , and block probability matrix

$$\tilde{B} = \begin{pmatrix} B_s & b_{sw} \mathbf{1} \\ b_{sw} \mathbf{1}' & 1 \end{pmatrix},$$

where the  $K \times K$  matrix  $B_s$ , as in Section 4.

Since  $\tilde{B}$  has rank  $K + 1$ , the same holds also for  $\tilde{\mathcal{L}}_{\tau}$ . We denote by  $\tilde{\mu}_{k,\tau}$ , for  $k = 1, \dots, n$ , to be the magnitude of the eigenvalues of  $\tilde{\mathcal{L}}_{\tau}$  arranged in decreasing order. Notice that  $\tilde{\mu}_{k,\tau} = 0$  for  $k > K + 1$ . Further, let  $\mathcal{V}_{\tau}$  be the  $n \times K$  eigenvector matrix of  $\tilde{\mathcal{L}}_{\tau}$ .

Lemma 11 shows that  $\tilde{\mu}_{2,\tau} = \dots = \tilde{\mu}_{K,\tau}$ , as well as provides explicit expression for these eigenvalues. Further, the lemma also characterizes the norm of the difference of the rows of  $\mathcal{V}_{\tau}$ . In the lemma below we denote by  $d_n^s = n^s p_n^s + (n - Kn^s)q_n + n^w b_{sw}$  and  $d_n^w = n^w + (n - n^w)b_{sw}$ . The quantities  $d_n^s$  and  $d_n^w$  provide the expected degrees of the nodes for an SBM drawn according to  $\tilde{B}$ .

**Lemma 11.** *The following holds:*

1. *The eigenvalue  $\tilde{\mu}_{1,\tau} = 1$ . Further, let  $\gamma_n = n^s(p_n^s - q_n)$ . Then*

$$\tilde{\mu}_{k,\tau} = \frac{\gamma_n}{d_n^s + \tau} \quad \text{for } k = 2, \dots, K \quad (47)$$

$$\tilde{\mu}_{K+1,\tau} = \frac{n^w(1 + \tau/n)}{d_n^w + \tau} - \frac{n^w(b_{sw} + \tau/n)}{d_n^s + \tau}. \quad (48)$$

2. The matrix  $\mathcal{V}_\tau$  has  $K+1$  distinct rows corresponding to the  $K+1$  clusters  $C_1, \dots, C_K$  and  $C^w$ . Denote these as  $\mathbf{cent}_{1,\tau}, \dots, \mathbf{cent}_{K,\tau}$  and  $\mathbf{cent}_\tau^w$ . Then  $1 \leq k' \neq k \leq K$

$$\|\mathbf{cent}_{k,\tau} - \mathbf{cent}_{k',\tau}\| = \sqrt{\frac{2}{n^s}}$$

for  $1 \leq k \leq K$ ,

$$\|\mathbf{cent}_{k,\tau} - \mathbf{cent}_\tau^w\| = \sqrt{\frac{1}{n^s}}$$

The above lemma is proved in Appendix C.3. Let  $\tilde{\mathcal{V}}_\tau$  be an  $n \times K$  matrix, with

$$\tilde{\mathcal{V}}_{i,\tau} = \mathbf{cent}_{k,\tau} \quad \text{for } i \in \tilde{C}_k.$$

Now  $\tilde{\mathcal{V}}_\tau$  has  $K$  distinct rows corresponding to the  $K$  clusters  $\tilde{C}_1, \dots, \tilde{C}_K$ . We denote these distinct rows as the population cluster centers. From Lemma 2, if

$$\|\mathbf{cent}_{k,\tau} - \mathbf{cent}_{k',\tau}\| \gtrsim (1/\delta) \|V_\tau - \tilde{\mathcal{V}}_\tau\| / \sqrt{n^s},$$

then  $\tilde{f} = O(\delta^2)$ . Since  $\|\mathbf{cent}_{k,\tau} - \mathbf{cent}_{k',\tau}\| \asymp 1/\sqrt{n^s}$  from Lemma 11, one gets that one needs to show that  $\|V_\tau - \tilde{\mathcal{V}}_\tau\| \lesssim \delta$ , with high probability, for some  $\delta$  that goes to zero for large  $n$ .

Now,

$$\begin{aligned} \|V_\tau - \tilde{\mathcal{V}}_\tau\| &\leq \|V_\tau - \mathcal{V}_\tau\| + \|\mathcal{V}_\tau - \tilde{\mathcal{V}}_\tau\| \\ &= \|V_\tau - \mathcal{V}_\tau\| + \sqrt{\frac{n^w}{n^s}} \end{aligned}$$

As  $n^w = O(1)$ , one needs to show that  $\|V_\tau - \mathcal{V}_\tau\|$  goes to zero with high probability. From Davis-Kahan theorem we get that

$$\begin{aligned} \|V_\tau - \mathcal{V}_\tau\| &\lesssim \frac{\|L_\tau - \tilde{\mathcal{L}}_\tau\|}{\tilde{\mu}_{K,\tau} - \tilde{\mu}_{K+1,\tau}} \\ &\lesssim \frac{\|L_\tau - \mathcal{L}_\tau\| + \|\mathcal{L}_\tau - \tilde{\mathcal{L}}_\tau\|}{\tilde{\mu}_{K,\tau} - \tilde{\mu}_{K+1,\tau}} \end{aligned} \quad (49)$$

The following lemma shows that for large  $\tau$ , the Laplacian matrix  $\mathcal{L}_\tau$  is close to the Laplacian matrix  $\tilde{\mathcal{L}}_\tau$  in spectral norm.

**Lemma 12.**

$$\|\mathcal{L}_\tau - \tilde{\mathcal{L}}_\tau\| \lesssim \frac{1}{1 + \tau/d_n^w}$$

The lemma is proved in Appendix C.2. Consequently, from Lemma 12 and Theorem 4 one gets from (49) that

$$\|V_\tau - \mathcal{V}_\tau\| \lesssim \frac{1}{(\tilde{\mu}_{K,\tau} - \tilde{\mu}_{K+1,\tau})} \left( \epsilon_{\tau,n} + \frac{1}{1 + \tau/d_n^w} \right) \quad (50)$$

Further, from Lemma 11 one gets that

$$\tilde{\mu}_{K,\tau} - \tilde{\mu}_{K+1,\tau} = \frac{n^s(p_n^s - q_n)}{d_n^s + \tau} - \left[ \frac{n^w(b_s + \tau/n)}{d_n^w + \tau} - \frac{n^w(b_{sw} + \tau/n)}{d_n^s + \tau} \right]$$

It is seen that  $(\tilde{\mu}_{K,\tau} - \tilde{\mu}_{K+1,\tau})\tau$  converges to

$$n^s(p_n^s - q_n) - [n^w(b_s - b_{sw}) + (n^w/n)(d_n^s - d_n^w)],$$

which is  $\gtrsim n^s(p_n^s - q_n)$  using  $n^w = O(1)$ .

Consequently, the right side of (50) converges to

$$\frac{\sqrt{d_n^s \log n + d_n^w}}{n^s(p_n^s - q_n)}$$

for large  $\tau$ . Now,  $d_n^s \asymp n p_n^s$  and  $d_n^w \asymp n b_{sw}$  (using  $n^w = O(1)$ ). Consequently, the numerator in the above is  $\lesssim \sqrt{n p_n^s}$  using Assumption (26). Consequently, under Assumption 24, one gets that  $\|V_\tau - \mathcal{V}_\tau\|$  goes to zero with high probability.

## C.2 Proof of Lemma 12

We bound the spectral norm of  $\mathcal{L}_\tau - \tilde{\mathcal{L}}_\tau$ . Here  $\tilde{\mathcal{L}}_\tau$  is as in Appendix C.1. Take  $\mathcal{L}_\tau = \mathcal{D}^{-1/2} (P + (\tau/n)J) \mathcal{D}^{-1/2}$  and  $\tilde{\mathcal{L}}_\tau = \tilde{\mathcal{D}}^{-1/2} (\tilde{P} + (\tau/n)J) \tilde{\mathcal{D}}^{-1/2}$ .

Notice that we ignore the subscript  $\tau$  in both  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$ . Here,  $\tilde{P} = Z\tilde{B}Z'$ , with  $\tilde{B}$  as in Subsection C.1.

As in the proof of Theorem 4, given in Appendix A.1, write

$$\mathcal{L}'_\tau = \tilde{\mathcal{D}}^{-1/2} (P + (\tau/n)J) \tilde{\mathcal{D}}^{-1/2}.$$

Then,

$$\|\mathcal{L}_\tau - \tilde{\mathcal{L}}_\tau\| \leq \|\mathcal{L}_\tau - \mathcal{L}'_\tau\| + \|\mathcal{L}'_\tau - \tilde{\mathcal{L}}_\tau\|. \quad (51)$$

Consequently, we prove that  $\mathcal{L}_\tau$  is close to  $\tilde{\mathcal{L}}_\tau$  by showing that both terms in the right side of (51) are small. We first bound  $\|\mathcal{L}_\tau - \mathcal{L}'_\tau\|$ . As in (38), write

$$\|\mathcal{L}_\tau - \mathcal{L}'_\tau\| \leq \|I - F\| (2 + \|I - F\|),$$

where as before  $F - I = \left( I + (\mathcal{D} - \tilde{\mathcal{D}}) \tilde{\mathcal{D}}^{-1} \right)^{1/2} - I$ . Here  $\mathcal{D} = \text{diag}(d_{1,\tau}, \dots, d_{n,\tau})$ , and  $\tilde{\mathcal{D}} = \text{diag}(\tilde{d}_{1,\tau}, \dots, \tilde{d}_{n,\tau})$ . Now,

$$\begin{aligned} \|(\mathcal{D} - \tilde{\mathcal{D}}) \tilde{\mathcal{D}}^{-1}\| &\leq \frac{|d_{i,\tau} - \tilde{d}_{i,\tau}|}{\tilde{d}_{i,\tau}} \\ &\lesssim \frac{d_n^w}{(d_n^w + \tau)}. \end{aligned}$$

Observe that we can assume that  $\|(\mathcal{D} - \tilde{\mathcal{D}}) \tilde{\mathcal{D}}^{-1}\| \leq 3/4$  for large  $\tau$ , so that

$$\left( 1 + \|(\mathcal{D} - \tilde{\mathcal{D}}) \tilde{\mathcal{D}}^{-1}\| \right)^{1/2} - 1 \leq \|(\mathcal{D} - \tilde{\mathcal{D}}) \tilde{\mathcal{D}}^{-1}\|,$$

and thus  $\|\mathcal{L}_\tau - \mathcal{L}'_\tau\| \lesssim d_n^w / (d_n^w + \tau)$ .

Next, we bound  $\|\mathcal{L}'_\tau - \tilde{\mathcal{L}}_\tau\|$ . Notice that  $\|\mathcal{L}'_\tau - \tilde{\mathcal{L}}_\tau\| \leq \|\tilde{\mathcal{G}}^{-1}\| \|(P - \tilde{P})\|$ . The quantity  $\|\tilde{\mathcal{G}}^{-1}\| \lesssim 1/(d_n^w + \tau)$ . Further, note that  $\|P - \tilde{P}\| \lesssim d_n^w$ , since  $P - \tilde{P}$  is a matrix with all entries negative and hence its spectral norm is at most the maximum of its row sums.

### C.3 Proof of Lemma 11

We investigate the eigenvalues of the  $K + 1$  community stochastic block model with block probability matrix

$$\tilde{B} = \begin{pmatrix} B_s & b_{sw}\mathbf{1} \\ b_{sw}\mathbf{1}' & b_w \end{pmatrix}$$

In our case  $b_w = 1$ . Denote the corresponding population Laplacian by  $\tilde{\mathcal{L}}$ . Recall that from Subsection A.2 the non-zero eigenvalues of  $\mathcal{L}$  are the same as that of

$$\tilde{B}_{eig} = (Z'RZ)^{1/2}B(Z'RZ)^{1/2}$$

Now,

$$Z'RZ = \text{diag}\left(\frac{n^s}{d_n^s}, \dots, \frac{n^s}{d_n^s}, \frac{n^w}{d_n^w}\right)$$

Consequently,

$$\tilde{B}_{eig} = \begin{pmatrix} \frac{n^s}{d_n^s} B_s & \left(\frac{n^s n^w}{d_n^s d_n^w}\right)^{1/2} b_{sw}\mathbf{1} \\ \left(\frac{n^s n^w}{d_n^s d_n^w}\right)^{1/2} b_{sw}\mathbf{1}' & \frac{n^w}{d_n^w} b_w \end{pmatrix},$$

One sees that

$$v_1 = (\sqrt{n^s d_n^s}, \dots, \sqrt{n^s d_n^s}, \sqrt{n^w d_n^w})'$$

is an eigenvector of  $\tilde{B}_{eig}$  with eigenvalue 1. Next, consider a vector  $v_2 = (v'_{21}, 0)'$ . Here  $v_{21}$  is a  $K \times 1$  dimensional vector that is orthogonal to the constant vector. We claim that  $v_2$  so defined is also an eigenvector of  $\tilde{B}_{eig}$ . To see this notice that

$$\tilde{B}_{eig} v_2 = \frac{n^s}{d_n^s} \begin{pmatrix} B_s v_{21} \\ 0 \end{pmatrix},$$

Here we use the fact that  $\mathbf{1}'v_{21} = 0$  as  $v_{21}$  is orthogonal to  $\mathbf{1}$ . Next, notice that

$$B_s = ((p_n^s - q_n)I + q_n \mathbf{1}\mathbf{1}')$$

Consequently,

$$B_s v_{21} = (p_n^s - q_n)v_{21}$$

The above implies that  $v_2$  is an eigenvector of  $\tilde{B}_{eig}$  with eigenvalue  $\lambda_1$  given by  $n^s(p_n^s - q_n)/d_n^s$ .

Notice that from the above construction one can get  $K - 1$  orthogonal eigenvectors  $v_k$ , for  $k = 2, \dots, K$ , such that the  $v_k$ 's are also orthogonal to  $v_1$ . Essentially, for  $k \geq 2$ , each  $v_k = (v'_{k1}, 0)'$ , where  $v'_{k1}\mathbf{1} = 0$ . There are  $K - 1$  orthogonal choices of the  $v_{k1}$ 's.

Given that 1 and  $\lambda_1$  are eigenvalues of  $\tilde{B}_{eig}$ , with the latter having multiplicity  $K - 1$ , the remaining eigenvalue is given by

$$\begin{aligned}\lambda_2 &= \text{trace}(\tilde{B}_{eig}) - 1 - (K - 1)\lambda_1 \\ &= \frac{n^s p_n^s}{d_n^s} + (K - 1)\frac{n^s}{d_n^s}q_n + \frac{n^w b_w}{d_n^w} - 1 \\ &= \frac{n^w b_w}{d_n^w} - \frac{n^w b_{sw}}{d_n^s}.\end{aligned}$$

The claim regarding the eigenvector corresponding to  $\lambda_2$  follows from seeing that this should be the case since it is orthogonal to eigenvectors  $v_1, \dots, v_K$  defined above.

## D Extending *DKest* to allow for degree heterogeneity

Here, we describe how we extend the *DKest* by substituting the estimate  $\hat{\mathcal{L}}_\tau$  in (30) with one assuming that the data is drawn from a degree corrected stochastic block model (D-SBM). As mentioned before, the D-SBM is a more appropriate model for modeling network datasets with extremely heterogeneous node degrees. The edge probability matrix takes the form

$$P = \Theta Z B Z' \Theta,$$

where  $\Theta = \text{diag}(\theta_1, \dots, \theta_n)$  models the heterogeneity in the degrees.

As before, assume that  $\hat{C}_{1,\tau}, \dots, \hat{C}_{K,\tau}$  be the cluster estimates obtained from running RSC- $\tau$  Algorithm. Let  $\hat{Z}$  be the corresponding  $n \times K$  cluster membership matrix. Denote

$$\hat{b}_{k_1, k_2} = \sum_{i \in \hat{C}_{k_1, \tau}, j \in \hat{C}_{k_2, \tau}} A_{ij}$$

and let  $\hat{B} = ((\hat{b}_{k_1, k_2}))$  be the  $K \times K$  with entries  $\hat{b}_{k_1, k_2}$ .

As in Karrer and Newman [14], we produce an estimate of the edge probability matrix  $P$  given by

$$\hat{P} = \hat{\Theta} \hat{Z} \hat{B} \hat{Z}' \hat{\Theta},$$

where  $\hat{\Theta} = \text{diag}(\hat{\theta}_1, \dots, \hat{\theta}_n)$ , with

$$\hat{\theta}_i = \frac{\hat{d}_i}{\sum_{k'=1}^K \hat{b}_{k, k'}}$$

for  $i \in \hat{C}_{k,\tau}$ . Recall that  $\hat{d}_i$  is the degree of node  $i$ . It is seen that with the above definition of  $\Theta$  the sum of the  $i$ -th row  $\hat{P}$  is simply  $\hat{d}_i$ .

The estimate  $\hat{\mathcal{L}}_\tau$  is taken as the population regularized Laplacian corresponding to the estimated edge probability matrix  $\hat{P}$ . In other words,

$$\hat{\mathcal{L}}_\tau = (D + \tau I)^{-1/2} \left( \hat{P} + \frac{\tau}{n} \mathbf{1}\mathbf{1}' \right) (D + \tau I)^{-1/2},$$

where recall that  $D$  is the diagonal matrix of degrees.

## References

- [1] Lada A Adamic and Natalie Glance. The political blogosphere and the 2004 us election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery*, pages 36–43. ACM, 2005.
- [2] A.A. Amini, A. Chen, P.J. Bickel, and E. Levina. Pseudo-likelihood methods for community detection in large sparse networks. *Ann. Statist.*, 41(4): 2097–2122, 2013.
- [3] Pranjal Awasthi and Or Sheffet. Improved spectral-norm bounds for clustering. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 37–49. Springer, 2012.
- [4] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- [5] Rajendra Bhatia. *Matrix analysis*, volume 169. Springer, 1997.
- [6] Peter J Bickel and Aiyu Chen. A nonparametric view of network models and newman–girvan and other modularities. *Proceedings of the National Academy of Sciences*, 106(50):21068–21073, 2009.
- [7] K. Chaudhuri, F. Chung, and A. Tsiatas. Spectral clustering of graphs with general degrees in the extended planted partition model. *Journal of Machine Learning Research*, 2012:1–23.
- [8] A. Chen, A. Amini, P. Bickel, and L. Levina. Fitting community models to large sparse networks. In *Joint Statistical Meetings, San Diego*, 2012.
- [9] Inderjit S Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In *Proc. seventh ACM SIGKDD inter. conf. on Know. disc. and data mining*, pages 269–274. ACM, 2001.
- [10] Donniell E Fishkind, Daniel L Sussman, Minh Tang, Joshua T Vogelstein, and Carey E Priebe. Consistent adjacency-spectral partitioning for the stochastic block model when the model parameters are unknown. *SIAM Journal on Matrix Analysis and Applications*, 34(1):23–39, 2013.

- [11] Lars Hagen and Andrew B Kahng. New spectral methods for ratio cut partitioning and clustering. *IEEE Trans. Computer-Aided Design*, 11(9):1074–1085, 1992.
- [12] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [13] A. Joseph and A.R. Barron. Fast sparse superposition codes have near exponential error probability for  $R < C$ . *IEEE. Trans. Inform. Theory*, to appear, 2013.
- [14] Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- [15] Amit Kumar and Ravindran Kannan. Clustering with spectral norm and the k-means algorithm. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 299–308. IEEE, 2010.
- [16] Tsz Chiu Kwok, Lap Chi Lau, Yin Tat Lee, Shayan Oveis Gharan, and Luca Trevisan. Improved cheeger’s inequality: Analysis of spectral partitioning algorithms through higher order spectral gap. *arXiv preprint arXiv:1301.5584*, 2013.
- [17] L. Mackey, M.I. Jordan, R.Y. Chen, B. Farrell, and J.A. Tropp. Matrix concentration inequalities via the method of exchangeable pairs. *arXiv preprint arXiv:1201.6002*, 2012.
- [18] Frank McSherry. Spectral partitioning of random graphs. In *Foundations of Computer Science, 2001. Proceedings. 42nd IEEE Symposium on*, pages 529–537. IEEE, 2001.
- [19] Mark EJ Newman and Michelle Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004.
- [20] Andrew Y Ng, Michael I Jordan, Yair Weiss, et al. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 2:849–856, 2002.
- [21] R.I. Oliveira. Concentration of the adjacency matrix and of the laplacian in random graphs with independent edges. *arXiv preprint arXiv:0911.0600*, 2009.
- [22] Tai Qin and Karl Rohe. Regularized spectral clustering under the degree-corrected stochastic blockmodel. *arXiv preprint arXiv:1309.4111*, 2013.
- [23] K. Rohe, S. Chatterjee, and B. Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *The Annals of Statistics*, 39(4):1878–1915, 2011.
- [24] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Trans. Pat. Analysis and Mach. Intel.*, 22(8):888–905, 2000.

- [25] Daniel L Sussman, Minh Tang, Donniell E Fishkind, and Carey E Priebe. A consistent adjacency spectral embedding for stochastic blockmodel graphs. *Journal of the American Statistical Association*, 107(499):1119–1128, 2012.
- [26] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [27] YY Yao. Information-theoretic measures for knowledge discovery and data mining. In *Entropy Measures, Maximum Entropy Principle and Emerging Applications*, pages 115–136. Springer, 2003.