

# Compressive sensing for cut improvement and local clustering.

Ming-Jun Lai\* and Daniel McKenzie †

**Abstract.** We show how one can phrase the cut improvement problem for graphs as a sparse recovery problem, whence one can use algorithms originally developed for use in compressive sensing (such as **SubspacePursuit** or **CoSaMP**) to solve it. We show that this approach to cut improvement is fast, both in theory and practice and moreover enjoys statistical guarantees of success when applied to graphs drawn from probabilistic models such as the Stochastic Block Model. Using this new cut improvement approach, which we call **ClusterPursuit**, as an algorithmic primitive we then propose new methods for local clustering and semi-supervised clustering, which enjoy similar guarantees of success and speed. Finally, we verify the promise of our approach with extensive numerical benchmarking.

**Key words.** Cluster Extraction, Local Clustering, Cut Improvement, Semi-Supervised Clustering, Community Detection, Compressive Sensing, Sparse Solution, Graph Laplacian.

**AMS subject classifications.** 68Q25, 68R10, 68U05, 94A12

**1. Introduction.** Finding clusters is a problem of primary interest when analyzing graphs. This is because vertices which are in the same cluster can reasonably be assumed to have some latent similarity. Thus, clustering can be used to find communities in social networks [24, 48, 53] or deduce political affiliation from a network of blogs [5]. Moreover, even data sets which are not presented as graphs can profitably be studied by first creating an auxiliary graph (eg. a  $K$ - or  $\epsilon$ -nearest-neighbors graph) and then applying graph clustering techniques. This has been successfully applied to image segmentation [43, 37], image classification [30] and natural language processing [19].

We shall informally think of a cluster as a subset of vertices,  $C \subset V$  with many edges between vertices in  $C$ , and few edges to the rest of the graph,  $C^c$ . See Figure 1 for a few examples. While some graphs may allow a neat partitioning into disjoint clusters (for example the OptDigits graph in Figure 1), for many graphs this is not the case. Some graphs may contain *background vertices*, that is, vertices which do not belong to any cluster (see the College Football graph in Figure 1). Alternatively, graphs may exhibit clusters at multiple scales (See the Senate Co-voting graph in Figure 1). In many cases, one has certain *a priori* information that could be used to improve clustering. For example in the OptDigits graph, we may know that some small subset,  $\Gamma \subset V$ , all represent images of ones. It is reasonable to assume that algorithms which incorporate this additional information (usually referred to as

---

\*mjlai@uga.edu. Department of Mathematics, University of Georgia, Athens, GA 30602. This research is partially supported by the National Science Foundation under the grant #DMS 1521537.

†mckenzie@math.ucla.edu. Department of Mathematics, University of California, Los Angeles, CA 155505. The financial assistance of the National Research Foundation of South Africa (NRF) towards this research is hereby acknowledged. Opinions expressed and conclusions arrived at, are those of the author and not necessarily to be attributed to the NRF. This research was conducted while this author was a graduate student at the University of Georgia and he gratefully acknowledges support and encouragement received from the Math Department of UGA.

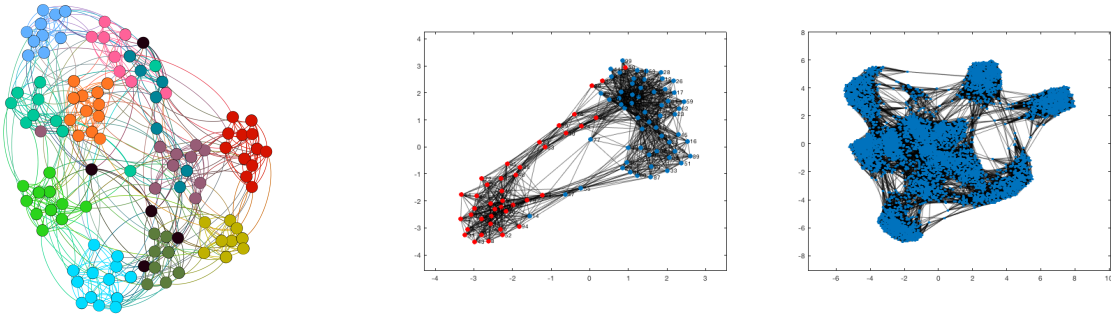


Figure 1: *Left*: the College Football graph of [24]. Vertices represents colleges fielding (American) football teams in the 2000 season. Vertices are connected if the respective teams played each other during the regular season. Clusters correspond to the various conferences in which teams play. Note that there are five schools, denoted in black, which are “independents” *ie* they are not affiliated with any conference. These can be thought of as background vertices. *Middle*: Senate co-voting for the 97th Congress, created using data from [32]. Vertices represent Senators and are connected if the respective Senators cast the same vote on a majority of bills. The two large clusters correspond to the two major American political parties. Notice how the blue cluster can be visually subdivided into two sub-clusters. *Right*: The OptDigits dataset consists of 5620 grayscale images of handwritten digits 0–9 of size  $8 \times 8$ . We discuss how to turn this into a graph in §10. Note that as there are ten digits, we expect this graph to have ten disjoint clusters.

35 semi-supervised algorithms) will perform better than ones which do not. With this in mind,  
 36 it is convenient to appeal to the following taxonomy of clustering algorithms:

- 37 1. *Global clustering algorithms* assign every vertex to one of  $k$  clusters, where the clusters  
 38 may or may not be disjoint. Algorithms for this problem may be unsupervised (for  
 39 example `SpectralClustering` [43, 40] or `GenLouvain` [18]) or semisupervised (for  
 40 example the auction dynamics approach of [30], or the regional force based methods  
 41 of [54]). This is appropriate for graphs such as the OptDigits graph of Figure 1, where  
 42 one expects a clear partition of the vertices into clusters.
- 43 2. *Local clustering algorithms*<sup>1</sup> take as input a small set of “seed vertices”,  $\Gamma \subset V$  and  
 44 return a good cluster containing  $\Gamma$ . Algorithms for local clustering are not confounded  
 45 by background vertices, as they are not required to assign them to a cluster. One  
 46 can further sub-divide local clustering algorithms into strongly and weakly local cluster-  
 47 ing algorithms. Strongly local algorithms, for `Nibble` [45, 46], `PPR-Grow` [1] or  
 48 `CapacityReleasingDiffusion` [52], are characterized by having run time proportional  
 49 to the size of the cluster found. This is advantageous when the cluster in question has  
 50 much fewer vertices than the graph as a whole. Weakly local algorithms are character-  
 51 ized as having run time proportional to the size of  $G$ . In practice they are frequently  
 52 faster than strongly local algorithms when finding large or moderately large clusters.

---

<sup>1</sup>Also known as cluster extraction algorithms in the statistics literature

We note that both kinds of local clustering algorithms may take as input a scale parameter, which dictates the size of the output cluster returned. This is useful when the graph at hand contains clusters at multiple scales, such as the Senate co-voting graph of Figure 1.

3. *Cut improvement algorithms* (cf. [1], [41], [50]) take as input a cut, or subset  $\Omega \subset V$ , which one can think of as an approximation to a cluster  $C$ , and refine it to produce a better approximation. Often cut improvement algorithms are run on the output of a local clustering algorithm to improve the quality of the output.

The central contribution of this paper is a new cut improvement algorithm which we call **ClusterPursuit**, that phrases the cut improvement problem as a sparse recovery problem. We pair this with a simple local clustering algorithm which we call Random Walk Thresholding or **RWThresh** to obtain a two-stage weakly local clustering algorithm that we shall refer to as **CP+RWT**. One can iterate this algorithm to find all clusters in a graph; we call this procedure iterated **CP+RWT** or **ICP+RWT**. After presenting some mathematical preliminaries and outlining the assumptions we place on generative models of graphs in §2, we derive the **ClusterPursuit** algorithm in §3 and prove that, given a cut  $\Omega$  satisfying  $|C_1 \Delta \Omega|/|C_1| = O(1)$  **ClusterPursuit** returns  $C_1^\#$  satisfying  $|C_1 \Delta C_1^\#|/|C_1| = o(1)$ . Here,  $C_1$  denotes the smallest cluster in the graph. In §4 we discuss the **RWThresh** algorithm, and show that given a small set of seed vertices,  $\Gamma \subset C_1$ , it is capable of finding an  $\Omega$  satisfying  $|C_1 \Delta \Omega|/|C_1| = O(1)$ . This leads naturally to guarantees of success for the two-stage local clustering algorithm **CP+RWT**, which we present in §5. In §6 we briefly discuss **ICP+RWT** while in §7 we show that **CP+RWT** and **ICP+RWT** enjoy a computational complexity of  $O(nd_{\max} \log(n))$  where  $d_{\max}$  is the largest vertex degree in the graph. In §8 we survey the literature and compare our work with relevant recent work in the area, while in §9 we show that a popular generative model of graphs with communities, namely the stochastic block model, satisfies the assumptions outlined in §2. Finally, we complement theoretical insight with experimental results in §10. In the interest of reproducibility, we make our code available at: [danielmckenzie.github.io](https://github.com/danielmckenzie).

## 2. Preliminaries.

**2.1. Graph Notation and Definitions.** We restrict our attention to finite, simple, undirected graphs  $G = (V, E)$ , possibly with non-negative edge weights. We identify the vertex set  $V$  with the integers  $[n] := \{1, \dots, n\}$  and denote an edge between vertices  $i$  and  $j$  as  $\{i, j\} \in E$ . The (possibly weighted) adjacency matrix of  $G$  will be denoted as  $A$ . By  $d_i$  we mean the degree of the  $i$ -th vertex, computed as  $d_i = \sum_j A_{ij}$ . For any  $S \subset V$  define  $\text{vol}(S) = \sum_{i \in S} d_i$ . For quantities such as  $d_i$  (and later  $\lambda_i$ ) that are indexed by  $i \in [n]$ , let  $d_{\max} := \max_i d_i$  and similarly  $d_{\min} := \min_i d_i$ . Denote by  $D$  the diagonal matrix whose  $(i, i)$  entry is  $d_i$ . By “cluster” we shall mean a subset of vertices,  $C \subset V$ , that is well-connected but sparsely connected to the rest of the graph. If a graph has clusters we shall refer to them as  $C_1, \dots, C_k$ . We define  $n_a := |C_a|$  and assume that the clusters are ordered by size, so that  $n_1 \leq n_2 \leq \dots \leq n_k$ . We reserve the letters  $a$  and  $b$  for indexing clusters, while  $i$  and  $j$  will index vertices.

94 **Definition 2.1 (Laplacians of graphs).** *The normalized, random walk Laplacian is defined as*  
 95  $L = I - D^{-1}A$ . *We shall simply refer to it as the Laplacian. The normalized, symmetric*  
 96 *Laplacian is:  $L^{sym} := I - D^{-1/2}AD^{-1/2}$ .*

97 Recall the following elementary result in spectral graph theory (see [49], for example, for  
 98 a proof):

99 **Theorem 2.2.** *Let  $C_1, \dots, C_k$  denote the connected components of a graph  $G$ . Then the*  
 100 *cluster indicator vectors  $\mathbf{1}_{C_1}, \dots, \mathbf{1}_{C_k}$  form a basis for the kernel of  $L$ .*

101 Suppose that  $G$  has clusters  $C_1, \dots, C_k$ . By definition, clusters have few edges between  
 102 them, and so it is useful to write  $G$  as the union of two edge-disjoint subgraphs, defined  
 103 as follows: let  $G^{\text{in}} = (V, E^{\text{in}})$  have only edges between vertices in the same cluster, while  
 104  $G^{\text{out}} = (V, E^{\text{out}})$  consist only of edges between vertices in different clusters. We emphasize  
 105 that this is a theoretical construction, as in practice we of course cannot ascertain whether two  
 106 vertices are in the same cluster without first solving the clustering problem, which is precisely  
 107 what we are trying to do. Denote by  $A^{\text{in}}$  and  $L^{\text{in}}$  (resp.  $A^{\text{out}}$  and  $L^{\text{out}}$ ) the adjacency  
 108 matrix and Laplacian of  $G^{\text{in}}$  (resp.  $G^{\text{out}}$ ). Similarly,  $d_i^{\text{in}}$  (resp.  $d_i^{\text{out}}$ ) shall denote the degree  
 109 of the vertex  $i$  in the graph  $G^{\text{in}}$  (resp.  $G^{\text{out}}$ ). For future reference we define the random  
 110 walk transition matrices  $P = AD^{-1}$  and  $N := D^{-1/2}AD^{-1/2}$ . We note that the spectra of  
 111  $P, N, A, L$  are related:

112 **Lemma 2.3.** *For any matrix  $B$  with real eigenvalues let  $\lambda_i(B)$  denote the  $i$ -th smallest*  
 113 *eigenvalue, counted with multiplicity. Then  $\lambda_i(L) = \lambda_i(L^{sym})$  while  $\lambda_{n-i}(N) = \lambda_{n-i}(P) =$   
 114  $1 - \lambda_i(L)$*

115 *Proof.* Observe that  $L = D^{-1/2}L^{sym}D^{1/2}$ , hence  $L$  and  $L^{sym}$  have the same spectrum.  
 116 Similarly  $P = D^{1/2}(I - L^{sym})D^{-1/2}$  hence  $P$  and  $N = I - L^{sym}$  have the same spectrum.  
 117 Thus if  $\lambda$  is the  $i$ -th smallest eigenvalue of  $L^{sym}$  it is the  $i$ -th largest (and hence the  $(n - i)$ -th  
 118 smallest) eigenvalue of  $I - L^{sym}$ . ■

119 For any  $S \subset V$ , we denote by  $G_S$  the induced sub-graph with vertices  $S$  and edges all  
 120  $\{i, j\} \in E$  with  $i, j \in S$ . By  $A_{G_S}$  (resp.  $L_{G_S}$ ) we mean the adjacency matrix (resp. Laplacian)  
 121 of the graph  $G_S$ . Note that  $L_{G_S}$  is not a submatrix of  $L$ ! For any  $S \subset [n]$  we define an *indicator*  
 122 *vector*  $\mathbf{1}_S \in \mathbb{R}^n$  by  $(\mathbf{1}_S)_i = 1$  if  $i \in S$  and  $(\mathbf{1}_S)_i = 0$  otherwise.  $|S|$  will always denote the  
 123 cardinality of  $S$ . For any matrix  $B$ , by  $B_S$  we mean the submatrix of  $B$  consisting of the  
 124 columns  $b_i$  for all  $i \in S$ .

125 **2.2. Compressive Sensing.** Recall for any  $\mathbf{x} \in \mathbb{R}^n$ ,  $\|\mathbf{x}\|_0 := |\text{supp}(\mathbf{x})| = |\{i : x_i \neq 0\}|$  is  
 126 the sparsity of  $\mathbf{x}$ . If  $\|\mathbf{x}\|_0 \ll n$  we say that  $\mathbf{x}$  is *sparse*. Candés, Donoho and their collaborators  
 127 in [20, 9] pioneered the study of compressive sensing, which offers theoretical analysis and  
 128 algorithmic tools for finding sparse solutions to linear systems  $\Phi\mathbf{x} = \mathbf{b}$ , for example by solving  
 129 the minimization problem:

$$130 \quad (2.1) \quad \text{argmin} \|\Phi\mathbf{x} - \mathbf{y}\|_2 \text{ subject to } \|\mathbf{x}\|_0 \leq s,$$

131 where  $\Phi \in \mathbb{R}^{m \times n}$  is referred to as the *sensing matrix*. Typically, it is assumed that  $m \leq n$   
 132 although this will not be the case in this paper. There are many algorithms available to solve

---

**Algorithm 2.1** SubspacePursuit, as presented in [17]

---

Input variables: measurement matrix  $\Phi$ , measurement vector  $\mathbf{y}$ , sparsity parameter  $s$  and number of iterations  $J$ .

Initialization:

- (1)  $S^{(0)} = \mathcal{L}_s(\Phi^\top \mathbf{y})$ .
- (2)  $\mathbf{x}^{(0)} = \arg \min_{\mathbf{z} \in \mathbb{R}^N} \{\|\mathbf{y} - \Phi \mathbf{z}\|_2 : \text{supp}(\mathbf{z}) \subset S^{(0)}\}$
- (3)  $\mathbf{r}^{(0)} = \mathbf{y} - \Phi \mathbf{x}^{(0)}$

**for**  $j = 1 : J$  **do**

- (1)  $\hat{S}^{(j)} = S^{(j-1)} \cup \mathcal{L}_s(\Phi^\top \mathbf{r}^{(j-1)})$
- (2)  $\mathbf{u} = \arg \min_{\mathbf{z} \in \mathbb{R}^N} \{\|\mathbf{y} - \Phi \mathbf{z}\|_2 : \text{supp}(\mathbf{z}) \subset \hat{S}^{(j)}\}$
- (3)  $S^{(j)} = \mathcal{L}_s(\mathbf{u})$  and  $\mathbf{x}^{(j)} = \mathcal{H}_s(\mathbf{u})$
- (4)  $\mathbf{r}^{(j)} = \mathbf{y} - \Phi \mathbf{x}^{(j)}$

**end for**

---

133 Problem (2.1), but the one we shall focus on is the **SubspacePursuit** algorithm introduced  
 134 in [17]. Here  $\mathcal{L}_s(\cdot)$  and  $\mathcal{H}_s(\cdot)$  are thresholding operators:

$$135 \quad \mathcal{L}_s(\mathbf{v}) := \{i \in [n] : v_i \text{ among } s \text{ largest-in-magnitude entries in } \mathbf{v}\}$$

$$136 \quad \mathcal{H}_s(\mathbf{v})_i := \begin{cases} v_i & \text{if } i \in \mathcal{L}_s(\mathbf{v}) \\ 0 & \text{otherwise.} \end{cases}$$

137

138 In quantifying whether (2.1) has a unique solution, the following constant is often used (see  
 139 [21])

140 **Definition 2.4.** The  $s$  Restricted Isometry Constant ( $s$ -RIC) of  $\Phi \in \mathbb{R}^{m \times n}$ , written  $\delta_s(\Phi)$ ,  
 141 is defined to be the smallest value of  $\delta > 0$  such that, for all  $\mathbf{x} \in \mathbb{R}^n$  with  $\|\mathbf{x}\|_0 \leq s$ , we have:

$$142 \quad (1 - \delta)\|\mathbf{x}\|_2^2 \leq \|\Phi \mathbf{x}\|_2^2 \leq (1 + \delta)\|\mathbf{x}\|_2^2.$$

143 If  $\delta_s(\Phi) < 1$  we often say that  $\Phi$  has the Restricted Isometry Property (RIP).

144 One of the reasons for the remarkable usefulness of compressive sensing is its robustness to  
 145 error, both additive (*i.e.* in  $\mathbf{y}$ ) and multiplicative (*i.e.* in  $\Phi$ ). More precisely, suppose that a  
 146 signal  $\hat{\mathbf{y}} = \hat{\Phi} \mathbf{x}^*$  is acquired, but that we do not know the sensing matrix  $\hat{\Phi}$  exactly. Instead,  
 147 we have access only to  $\Phi = \hat{\Phi} + M$ , for some small perturbation  $M$ . Suppose further that  
 148 there is some noise in the measurement process, so that the signal we actually receive is  
 149  $\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e}$ . Can one hope to approximate a sparse vector  $\mathbf{x}^*$  from  $\mathbf{y}$ , given only  $\Phi$ ? This  
 150 question is answered in the affirmative way by several authors, starting with the work of [29].  
 151 For **SubspacePursuit**, we have the following result (cf. [33]):

152 **Theorem 2.5.** Let  $\mathbf{x}^*$ ,  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ ,  $\Phi$  and  $\hat{\Phi}$  be as above and suppose that  $\|\mathbf{x}^*\|_0 \leq s$ . For any  
 153  $t \in [n]$ , let  $\delta_t := \delta_t(\Phi)$ . Define the following constants:

$$154 \quad \epsilon_{\mathbf{y}} := \|\mathbf{e}\|_2 / \|\hat{\mathbf{y}}\|_2 \text{ and } \epsilon_{\Phi}^s = \|M\|_2^{(s)} / \|\hat{\Phi}\|_2^{(s)}$$

155 where for any matrix  $B$ ,  $\|B\|_2^{(s)} := \max\{\|B_S\|_2 : S \subset [n] \text{ and } |S| = s\}$ . Define further:

$$156 \quad \rho = \frac{\sqrt{2\delta_{3s}^2(1 + \delta_{3s}^2)}}{1 - \delta_{3s}^2} \quad \text{and} \quad \tau = \frac{(\sqrt{2} + 2)\delta_{3s}}{\sqrt{1 - \delta_{3s}^2}}(1 - \delta_{3s})(1 - \rho) + \frac{2\sqrt{2} + 1}{(1 - \delta_{3s})(1 - \rho)}$$

157 Assume  $\delta_{3s} \leq 0.4859$  and let  $\mathbf{x}^{(m)}$  be the output of `SubspacePursuit` applied to Problem (2.1)  
158 after  $m$  iterations. Then:

$$159 \quad \frac{\|\mathbf{x}^* - \mathbf{x}^{(m)}\|_2}{\|\mathbf{x}^*\|_2} \leq \rho^m + \tau \frac{\sqrt{1 + \delta_s}}{1 - \epsilon_{\Phi}^s} (\epsilon_{\Phi}^s + \epsilon_{\mathbf{y}}).$$

160 *Proof.* This is Corollary 1 in [33]. Note that our convention on hats is different to theirs  
161 — our  $\Phi$  is their  $\hat{\Phi}$ , hence our  $\rho$  is their  $\hat{\rho}$  and so on. ■

162 Next it is easy to obtain bounds on the quantity  $\|B\|_2^{(s)} := \max_{\substack{S \subset [n] \\ |S|=s}} \|B_S\|_2$ :

163 **Lemma 2.6.** For any matrix  $B$  and any  $2 \leq s \leq n$  we have that  $\sigma_{s-1}(B) \leq \|B\|_2^{(s)} \leq$   
164  $\sigma_{\max}(B) = \|B\|_2$ , where  $\sigma_j(B)$  denotes the  $j$ -th smallest singular value of  $B$ .

165 *Proof.* Observe that, for any matrix  $B$ ,

$$166 \quad \|B\|_2^{(s)} = \max_{\substack{S \subset [n] \\ |S|=s}} \|B_S\|_2 = \max_{\substack{S \subset [n] \\ |S|=s}} \sigma_{\max}(B_S),$$

167 where  $\sigma_{\max}(B_S)$  denotes the maximum singular value of  $B_S$ . Because  $\sigma_{\max}(B_S) = \sigma_s(B_S)$ , by  
168 the interlacing theorem for singular values (cf. [47])  $\sigma_{s-1}(B) \leq \sigma_{\max}(B_S) \leq \sigma_{\max}(B)$ . ■

169 **2.3. The Data Model.** For conceptual clarity, we shall take an asymptotic viewpoint,  
170 and consider graphs  $G \in \mathcal{G}_n$  as  $n \rightarrow \infty$ . Note that the graphs under consideration may be  
171 weighted or unweighted. We say that a graph property  $P$  holds *almost surely* for  $\mathcal{G}_n$  if the  
172 probability of a  $G$  drawn from  $\mathcal{G}_n$  *not having*  $P$  is  $o(1)$ .

173 **Assumptions 2.7.** Suppose that there exist  $\epsilon_i = o(1)$  as  $n \rightarrow \infty$  for  $i = 1, 2, 3$  such that for  
174 all  $G \in \mathcal{G}_n$ :

- 175 (A1)  $V = C_1 \cup \dots \cup C_k$  where the  $C_a$  are disjoint, planted clusters and  $k$  is  $O(1)$  as  $n \rightarrow \infty$ .
- 176 (A2) For all  $a \in [k]$  we have that  $\lambda_2(L_{G_{C_a}}) \geq 1 - \epsilon_1$  and  $\lambda_{n_a}(L_{G_{C_a}}) \leq 1 + \epsilon_1$  almost surely.
- 177 (A3) Letting  $r_i := d_i^{\text{out}}/d_i^{\text{in}}$ ,  $r_i \leq \epsilon_2$  for all  $i \in [n]$  almost surely.
- 178 (A4) If  $d_{\text{av}}^{\text{in}} := \mathbb{E}[d_i^{\text{in}}]$  then  $d_{\max}^{\text{in}} \leq (1 + \epsilon_3)d_{\text{av}}^{\text{in}}$  and  $d_{\min}^{\text{in}} \geq (1 - \epsilon_3)d_{\text{av}}^{\text{in}}$  almost surely.

179 Note that we can think of (A1)–(A4) as “regularity” requirements for graphs; as they  
180 insist that degrees do not vary too wildly, and that the eigenvalues are well behaved. In  
181 §9 we verify that a common model of unweighted graphs with clusters—the stochastic block  
182 model—satisfies these assumptions, so they are certainly not too restrictive. It seems probable  
183 (and indeed supported by the numerical evidence of §10) that reasonable models of random  
184 weighted graphs satisfy these properties too, although we leave this for future work.

185 **3. The ClusterPursuit Algorithm.** The motivation for our algorithm is the following  
 186 observation. Suppose for a moment that one had access to  $L^{\text{in}}$ . Suppose further that one is  
 187 given a cut  $\Omega$  “near” a cluster of interest,  $C_a$ , which we shall take quantitatively to mean that  
 188  $|C_a \Delta \Omega| = \epsilon|C_a|$ , where  $\Delta$  denotes the symmetric difference, i.e  $C \Delta \Omega = (C \setminus \Omega) \cup (\Omega \setminus C)$  and  
 189  $\epsilon \in (0, 1)$ . Letting  $U = C_a \setminus \Omega$  and  $W = \Omega \setminus C_a$  one observes that:

$$\begin{aligned} 190 \quad & \mathbf{1}_\Omega = \mathbf{1}_{C_a} + \mathbf{1}_W - \mathbf{1}_U \\ 191 \quad & \implies L^{\text{in}}\mathbf{1}_\Omega = L^{\text{in}}\mathbf{1}_{C_a} + L^{\text{in}}(\mathbf{1}_W - \mathbf{1}_U) \\ 192 \quad & \implies L^{\text{in}}\mathbf{1}_\Omega = 0 + L^{\text{in}}(\mathbf{1}_W - \mathbf{1}_U) \quad (\text{by Theorem 2.2}) \\ 193 \quad & \implies \mathbf{y}^{\text{in}} = L^{\text{in}}(\mathbf{1}_W - \mathbf{1}_U) \quad (\text{if } \mathbf{y}^{\text{in}} := L^{\text{in}}\mathbf{1}_\Omega) \end{aligned}$$

195 Solving the linear system  $\mathbf{y}^{\text{in}} = L^{\text{in}}\mathbf{x}$  is unlikely to yield  $\mathbf{x} = \mathbf{1}_W - \mathbf{1}_U$ , as  $L^{\text{in}}$  has a large  
 196 kernel (Theorem 2.2). However, Theorem 3.2 will show that one may recover  $\mathbf{1}_W - \mathbf{1}_U$  as the  
 197 solution to the sparse recovery problem:

$$198 \quad (3.1) \quad \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{ \|L^{\text{in}}\mathbf{x} - \mathbf{y}^{\text{in}}\|_2 : \|\mathbf{x}\|_0 \leq s \}$$

199 where  $s \approx |C_a \Delta \Omega|$ . Of course, one will not in practice have access to  $L^{\text{in}}$ , only  $L$ . Thus one  
 200 needs to consider a perturbed version of (3.1):

$$201 \quad (3.2) \quad \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{ \|L\mathbf{x} - \mathbf{y}\|_2 : \|\mathbf{x}\|_0 \leq s \}$$

202 where  $\mathbf{y} = L\mathbf{1}_\Omega$ . Theorem 3.4 will show that the solution  $\mathbf{x}^\#$  to the minimization problem  
 203 (3.2) found by `SubspacePursuit` is a good enough approximation to  $\mathbf{1}_U - \mathbf{1}_W$ , hence one  
 204 may infer  $U$  and  $W$  from the signed support of  $\mathbf{x}^\#$ . Clearly, if one knows  $\Omega, U$  and  $W$  one  
 205 may reconstruct  $C_a$  as  $C_a = (\Omega \setminus W) \cup U$ . This is the essence of `ClusterPursuit`, which we  
 206 present as Algorithm 3.1.

---

#### Algorithm 3.1 ClusterPursuit

---

**Input:** Adjacency matrix  $A$ , initial cut  $\Omega$ , estimate  $s \approx |\Omega \Delta C_a|$  and  $R \in [0, 1)$ .

(1) Compute  $L = I - D^{-1}A$  and  $\mathbf{y} = L\mathbf{1}_\Omega$ .

(2) Let  $\mathbf{x}^\#$  be the solution to

$$(3.3) \quad \operatorname{argmin} \{ \|L\mathbf{x} - \mathbf{y}\|_2 : \|\mathbf{x}\|_0 \leq s \}$$

obtained after  $m = O(\log(n))$  iterations of `SubspacePursuit`.

(3) Let  $U^\# = \{i : x_i^\# < -R\}$  and  $W^\# = \{i : x_i^\# > R\}$ .

**Output:**  $C_a^\# = (\Omega \setminus W^\#) \cup U^\#$ .

---

207 *Remark 3.1.* `ClusterPursuit` requires as an input an estimate of  $|\Omega \Delta C_a|$ , which might  
 208 not always be available. This is less of an issue than it might first appear as:

- 209 1. Theorem 3.4 will show that as long as  $|\Omega \triangle C_a| \leq s \leq 0.13|C_a|$  `ClusterPursuit` works  
 210 well.  
 211 2. If no knowledge of  $|\Omega \triangle C_a|$  is available, one may run `ClusterPursuit` for various  
 212 values of  $s$  and keep the returned cluster with lowest conductance.  
 213 3. Alternatively, one could consider the Lasso form of problem (3.3):

$$214 \quad (3.4) \quad \operatorname{argmin} \{ \|L\mathbf{x} - \mathbf{y}\|_2 + \lambda \|\mathbf{x}\|_1 \} = \operatorname{argmin} \{ \|L\mathbf{x} - \mathbf{y}\|_2 + \lambda \|\mathbf{x}\|_0 \}$$

215 as the sparse solution is the cluster indicator  $\mathbf{1}_U - \mathbf{1}_W$  which satisfies  $\|\mathbf{x}\|_0 = \|\mathbf{x}\|_1$ .  
 216 We do not analyze this further here.

217 **Theorem 3.2.**  $\mathbf{1}_W - \mathbf{1}_U$  is the unique solution to Problem (3.1), for any graph  $G$  with  
 218 clusters  $C_1, \dots, C_k$ , as long as  $|C_a \triangle \Omega| \leq s < n_1/2$ .

219 *Proof.* One can easily verify that  $\mathbf{1}_W - \mathbf{1}_U$  is a solution to (3.1), thus it remains to show  
 220 that it is the unique one. So, suppose that  $\mathbf{v}$  satisfies  $L^{\text{in}}\mathbf{v} = \mathbf{y}^{\text{in}}$  and that  $\mathbf{v} \neq \mathbf{1}_W - \mathbf{1}_U$ .  
 221 Because  $\mathbf{y}^{\text{in}} = L^{\text{in}}\mathbf{1}_\Omega$ :

$$222 \quad L^{\text{in}}\mathbf{v} - L^{\text{in}}\mathbf{1}_\Omega = 0 \implies \mathbf{v} - \mathbf{1}_\Omega \in \ker(L^{\text{in}}) \implies \mathbf{v} - \mathbf{1}_\Omega = \sum_{b=1}^k \alpha_b \mathbf{1}_{C_b} \quad (\text{by Theorem 2.2})$$

$$223 \implies \mathbf{v} = \sum_{b=1}^k \alpha_b \mathbf{1}_{C_b \setminus \Omega} + \sum_{b=1}^k (\alpha_b + 1) \mathbf{1}_{C_b \cap \Omega}$$

$$224$$

225 Now if  $\alpha_a = -1$  and  $\alpha_b = 0$  for all  $b \neq a$  then  $\mathbf{v} = \mathbf{1}_W - \mathbf{1}_U$ , which we are assuming is not  
 226 the case. Hence either  $\alpha_a \neq -1$ , in which case  $\|\mathbf{v}\|_0 \geq |C_a \cap \Omega| \geq |C_a| - |C_a \triangle \Omega| > n_a/2$ , or  
 227  $\alpha_b \neq 0$  for  $b \neq a$  in which case  $\|\mathbf{v}\|_0 \geq |C_b \setminus \Omega| \geq |C_b| - |C_a \triangle \Omega| > n_b/2$  as we are assuming  
 228 that  $|C_a \triangle \Omega| < n_1/2$  and  $n_1 = \min_b n_b$ . By assumption,  $s < n_1/2$  hence in either case  $\mathbf{v}$  is  
 229 infeasible for Problem (3.1), as it does not satisfy the constraint  $\|\mathbf{v}\|_0 \leq s$ .  $\blacksquare$

230 Henceforth, we shall focus on recovering the smallest cluster,  $C_1$ . We do this to avoid a  
 231 technical complication in the estimation of  $\delta_{\gamma n_a}(L)$  for  $a > 1$  (see Theorem 3.3 and Remark  
 232 A.3). We note that as long as  $n_a \approx n_1$  this is not really an issue, and the proof of Theorem  
 233 3.4 will extend to this case, albeit with a tighter bound on  $\epsilon$ .

234  
 235 Let us now quantify the size of the perturbation in moving from (3.1) to (3.2). Define  
 236  $M := L - L^{\text{in}}$  and  $\mathbf{e} := \mathbf{y} - \mathbf{y}^{\text{in}}$ . Recall from Theorem 2.5, that the three key parameters in  
 237 perturbed compressive sensing are the restricted isometry constant of  $L$  and:

$$238 \quad (3.5) \quad \epsilon_{\mathbf{y}} = \frac{\|\mathbf{e}\|_2}{\|\mathbf{y}^{\text{in}}\|_2} \quad \text{and} \quad \epsilon_L^s = \frac{\|M\|_2^{(s)}}{\|L^{\text{in}}\|_2^{(s)}}$$

239 as well as two secondary quantities,  $\rho$  and  $\tau$ . We prove the following:

240 **Theorem 3.3.** Suppose that  $\mathcal{G}_n$  satisfies (A1)–(A4) and that  $|\Omega \triangle C_1| \leq 0.13n_1$ . Then for  
 241 any  $\gamma \in (0, 1)$  the following hold almost surely:

- 242 1.  $\epsilon_{\mathbf{y}} = o(1)$  and  $\epsilon_L^{\gamma n_1} = o(1)$ .



- 243 2.  $\delta_{\gamma n_1}(L) \leq \gamma + o(1)$ .  
 244 3. If  $\delta_{3s}(L) \leq 0.45$  then  $\rho \leq 0.8751$  and  $\tau \leq 55.8490$  for any  $s \in (0, n_1/3)$ .

245 *Proof.* Part (3) follows by direct computation. For parts (1) and (2) see Appendix A. ■

246 We now prove the main result of this section:

**Theorem 3.4.** *Suppose that  $A$  is the adjacency matrix of  $G \sim \mathcal{G}_n$  satisfying assumptions (A1)–(A4), and that  $\Omega$  satisfies  $|C_1 \triangle \Omega| = \epsilon n_1$  with  $\epsilon \leq 0.13$ . If  $C^\#$  is the output of ClusterPursuit when given inputs  $A, \Omega, \epsilon n_1 \leq s \leq 0.13n_1$  and  $R = 0.5$  then:*

$$\frac{|C_1 \triangle C_1^\#|}{|C_1|} = o(1) \quad \text{almost surely.}$$

247 *Remark 3.5.*  $s \leq 0.13n_1$  is a conservative upper bound on  $s$  for which the guarantees from  
 248 §2.2 will hold. If one has no further information on  $|C_1 \triangle \Omega|$  we recommend using this as  
 249 the default value of  $s$ . Empirically (see §10.1) we still observe excellent performance when  
 250  $s > 0.13n_1$ .

251 *Proof.* Recall  $\mathbf{x}^\#$  is the solution obtained by  $m = O(\log(n))$  iterations of SubspacePursuit  
 252 on Problem 3.2, which we are regarding as a perturbation of Problem 3.1. Clearly,  $0.13n_1 <$   
 253  $n_1/2$ , hence by Theorem 3.2  $\mathbf{1}_W - \mathbf{1}_U$  is the unique solution to (3.1). By Theorem 3.3 part (2)  
 254 we get  $\delta_s(L) \leq 0.13 + o(1) < 0.15$  almost surely, for large enough  $n_1$ . Similarly  $\delta_{3s}(L) \leq 0.45$ ,  
 255 again almost surely for  $n_1$  large enough. It follows from Theorem 3.3 part (3) that  $\rho \leq 0.8751$   
 256 and  $\tau \leq 55.8490$ . We now appeal to Theorem 2.5 to obtain:

$$257 \frac{\|(\mathbf{1}_U - \mathbf{1}_W) - \mathbf{x}^\#\|_2}{\|\mathbf{1}_U - \mathbf{1}_W\|_2} \leq \rho^m + \tau \frac{\sqrt{1 + \delta_s}}{1 - \epsilon_\Phi^s} (\epsilon_\Phi^s + \epsilon_Y).$$

258 The second term on the right-hand side is  $o(1)$  by Theorem 3.3. As long as  $m \geq \log_\rho(1/n) =$   
 259  $O(\log(n))$ , we obtain that  $\rho^m = 1/n = o(1)$  too. Thus:

$$260 (3.6) \quad \frac{\|(\mathbf{1}_U - \mathbf{1}_W) - \mathbf{x}^\#\|_2}{\|\mathbf{1}_U - \mathbf{1}_W\|_2} \leq o(1) \implies \|(\mathbf{1}_U - \mathbf{1}_W) - \mathbf{x}^\#\|_2 \leq o(\|\mathbf{1}_U - \mathbf{1}_W\|_2) = o(\sqrt{n_1})$$

261 As  $\|\mathbf{1}_U - \mathbf{1}_W\|_2 = \sqrt{|U| + |W|} = \sqrt{\epsilon n_1}$ . In Lemma 3.6 below we show that, because  $\mathbf{1}_U - \mathbf{1}_W$  is  
 262 a difference of binary vectors, equation (3.6) implies that  $|U \triangle U^\#| = o(n_1)$  and  $|W \triangle W^\#| =$   
 263  $o(n_1)$ , and hence  $|C_1 \triangle C_1^\#| = o(n_1)$ , as required. ■

**Lemma 3.6.** *Consider disjoint  $T_1, T_2 \subset [n]$  and any  $\mathbf{v} \in \mathbb{R}^n$ . Define  $T_1^\# = \{i : v_i > 0.5\}$   
 and  $T_2^\# = \{i : v_i < -0.5\}$ . If  $\|(\mathbf{1}_{T_1} - \mathbf{1}_{T_2}) - \mathbf{v}\|_2 \leq D$  then:*

$$|T_1 \triangle T_1^\#| + |T_2 \triangle T_2^\#| \leq 4D^2.$$

*Proof.* Let  $T_3 := [n] \setminus (T_1 \cup T_2)$  and write  $\mathbf{v} = \mathbf{v}^{(1)} + \mathbf{v}^{(2)} + \mathbf{v}^{(3)}$  where  $\mathbf{v}^{(i)}$  denotes the  
 part of  $\mathbf{v}$  supported on  $T_i$ . Observe that:

$$D^2 \geq \|(\mathbf{1}_{T_1} - \mathbf{1}_{T_2}) - \mathbf{v}\|_2^2 = \|\mathbf{1}_{T_1} - \mathbf{v}^{(1)}\|_2^2 + \|\mathbf{1}_{T_2} - \mathbf{v}^{(2)}\|_2^2 + \|\mathbf{v}^{(3)}\|_2^2$$

One can easily verify that:

$$\|\mathbf{1}_{T_1} - \mathbf{v}^{(1)}\|_2^2 = \|\mathbf{1}_{T_1 \cap T_1^\#} - \mathbf{v}^{(1)}|_{T_1 \cap T_1^\#}\|_2^2 + \|\mathbf{1}_{T_1 \setminus T_1^\#} - \mathbf{v}^{(1)}|_{T_1 \setminus T_1^\#}\|_2^2 \geq (0.5)^2 |T_1 \setminus T_1^\#|$$

Similarly,  $\|-\mathbf{1}_{T_2} + \mathbf{v}^{(2)}\|_2^2 \geq (0.5)^2 |T_2 \setminus T_2^\#|$ , and:

$$\|\mathbf{v}^{(3)}\|_2^2 \geq \|\mathbf{v}^{(3)}|_{T_1^\# \setminus T_1}\|_2^2 + \|\mathbf{v}^{(3)}|_{T_2^\# \setminus T_2}\|_2^2 \geq (0.5)^2 |T_1^\# \setminus T_1| + (0.5)^2 |T_2^\# \setminus T_2|$$

Putting this all together we get that:

$$D^2 \geq (0.5)^2 \left( |T_1 \setminus T_1^\#| + |T_2 \setminus T_2^\#| + |T_1^\# \setminus T_1| + |T_2^\# \setminus T_2| \right) = 0.25 \left( |T_1 \triangle T_1^\#| + |T_2 \triangle T_2^\#| \right)$$

264 **4. The RWThresh algorithm.** Here, we introduce a simple, diffusion-based local cluster-  
 265 ing algorithm which we call **RWThresh** (see Algorithm 4.1). We note that **RWThresh** is some-  
 266 what similar to other more sophisticated diffusion-based local clustering algorithms, such as  
 267 **PPR-Grow**, **HK-Grow** and **CapacityReleasingDiffusion**. We do not claim that **RWThresh**  
 268 outperforms similar existing algorithms; its main utility lies in the fact that it reliably (and  
 269 provably) produces approximate cuts,  $\Omega$ , that are of a high enough quality to be used as an  
 270 initialization for **ClusterPursuit**.

---

#### Algorithm 4.1 RWThresh

---

**Input:** Adjacency matrix  $A$ , a thresholding parameter  $\epsilon \in (0, 1)$ , seed vertices  $\Gamma \subset C_1, \hat{n}_1 \approx n_1$  and depth of random walk  $t$ .

(1) Compute  $P = AD^{-1}$  and let  $\mathbf{v}^{(0)} = D\mathbf{1}_\Gamma$ .

(2) Compute  $\mathbf{v}^{(t)} = P^t \mathbf{v}^{(0)}$

(3) Define  $\Omega = \tilde{\mathcal{L}}_{(1+\epsilon)\hat{n}_1}(\mathbf{v}^{(t)})$ .

**Output:**  $\Omega = \Omega \cup \Gamma$ .

---

271 Here  $\tilde{\mathcal{L}}_t(\cdot)$  is a thresholding operator, similar to  $\mathcal{L}_t(\cdot)$ , but that returns the indices of the  $t$   
 272 largest, not largest-in-magnitude, components of a vector. To motivate **RWThresh** we observe  
 273 the following. If  $\mathcal{G}_n$  satisfies assumptions (A1)–(A4) then:

- 274 1.  $G_{C_1}$  is sufficiently densely connected that after  $t$  steps the random walk has a fairly  
 275 large probability of visiting every  $i \in C_1$ .
- 276 2.  $C_1$  is sufficiently weakly connected to  $V \setminus C_1$  that the probability of the random walk  
 277 leaving  $C_1$  after  $t$  steps is fairly small.

278 Hence Algorithm, 4.1 which runs a short random walk starting on  $\Gamma$  and takes  $\Omega$  to be the  
 279 set of vertices most likely to be visited, should produce an  $\Omega$  which is close to our intuitive  
 280 notion of a good cluster. Let us quantify this as Theorem 4.1.

281 **Theorem 4.1.** *Let  $G \sim \mathcal{G}_n$  satisfy Assumptions (A1)–(A4) and let  $A$  denote the adjacency  
 282 matrix of  $G$ . Let  $\Omega$  denote the output of **RWThresh** with inputs  $A$ , any  $\epsilon \in (0, 1)$ , any  $t = O(1)$ ,  
 283  $\hat{n}_1 = n_1$  and  $\Gamma \subset C_1$  with  $|\Gamma| = g\epsilon_3^{2t-1}n_1$  for any constant  $g \in (0, 1)$ , where  $\epsilon_3$  is as in  
 284 Assumption (A4)). Then  $|\Omega \triangle C_1| \leq (\epsilon + o(1))n_1$  almost surely.*

285 *Proof.* The proof is left to Appendix B. ■

286 We note that there are many local clustering algorithms, for example the PPR-Grow and  
 287 CapacityReleasingDiffusion algorithms discussed in §8, that require only  $|\Gamma| = O(1)$ .  
 288 However, these algorithms tend to return small clusters, typically of size  $|C| = O(1)$ . If  
 289  $\epsilon_3 = O(1/\log(n))$ , as it is in the numerical experiments of §10.1, then Theorem 4.1 requires  
 290 that  $|\Gamma| = O(n_1/\text{polylog}(n_1))$ , which seems to be a reasonable assumption when finding a  
 291 cluster of size  $O(n)$ . In practice, we find it suffices to take  $|\Gamma| = 0.01n_1$  or  $|\Gamma| = 0.02n_1$ .

292 **5. Using ClusterPursuit for local clustering.** As mentioned earlier, using RWThresh to  
 293 quickly generate a rough approximation to  $C_1$ , namely  $\Omega$ , and then using ClusterPursuit to  
 294 then refine this cut leads to a (weakly) local clustering algorithm. Here we verify that this  
 295 approach, presented below as algorithm 5.1, works well for our model of graph.

---

**Algorithm 5.1** CP+RWT
 

---

**Input:** Adjacency matrix  $A$  and seed vertices  $\Gamma \subset C_1$ . Parameters  $\epsilon \in (0, 0.13)$ ,  $s \approx \epsilon n_1$ ,  
 $R \in [0, 1)$ ,  $\hat{n}_1 \approx n_1$ ,  $t \in \mathbb{Z}_+$

(1) Let  $\Omega = \text{RWThresh}(A, \epsilon, \Gamma, \hat{n}_1, t)$

(2) Let  $C_1^\# = \text{ClusterPursuit}(A, s, R)$

**Output:**  $C_1^\#$

---

**Theorem 5.1.** *Let  $G \sim \mathcal{G}_n$  satisfy Assumptions (A1)–(A4) and let  $A$  denote the adjacency matrix of  $G$ . Let  $C_1^\#$  denote the output of CP+RWT with inputs  $A$ ,  $\epsilon \in (0, 0.13)$ ,  $R = 0.5$ ,  $\hat{n}_1 = n_1$ , any  $t = O(1)$ , any  $s$  satisfying  $\epsilon < s \leq 0.13n_1$  and  $\Gamma \subset C_1$  with  $|\Gamma| = g\epsilon_3^{2t-1}n_1$  for any constant  $g \in (0, 1)$ , where  $\epsilon_3$  is as in Assumption (A4). Then:*

$$\frac{|C_1 \Delta C_1^\#|}{|C_1|} = o(1)$$

296 *almost surely, for large enough  $n_1$ .*

297 *Proof.* By Theorem 4.1, the call to RWThresh in Step (1) of CP+RWT almost surely returns  
 298 an  $\Omega$  satisfying  $|\Omega \Delta C_1| \leq (\epsilon + o(1))n_1$  for input parameters with the given values. For large  
 299 enough  $n_1$ , we have that  $(\epsilon + o(1))n_1 \leq s \leq 0.13n_1$ , hence the call to ClusterPursuit in Step  
 300 (2) of CP+RWT returns  $C_1^\#$  with  $|C_1 \Delta C_1^\#|/|C_1| = o(1)$  by Theorem 3.4, again almost surely. ■

301 **Remark 5.2.** In practice (see §10) we find it generally suffices to take  $t = 3$ . If  $C_1$  is  
 302 densely connected, one might consider a smaller value of  $t$ , and conversely one might choose  
 303 a larger value (say  $t = 5$ ) if  $C_1$  is sparsely connected.

304 **6. Using ClusterPursuit for semi-supervised clustering.** In the (global) semi-supervised  
 305 clustering problem, one is given a small set of seed vertices  $\Gamma_a \subset C_a$  in each cluster, usu-  
 306 ally referred to in this context as “labeled data”. The goal here is to find a partition into  
 307 disjoint sets:  $V = C_1^\# \cup C_2^\# \cup \dots \cup C_k^\#$  that closely resembles the ground truth partition  
 308  $V = C_1 \cup C_2 \cup \dots \cup C_k$ . An iterated version of CP+RWT, which we call ICP+RWT, can be used  
 309 to solve this problem. ICP+RWT is presented as algorithm 6.1. Note that in the second line  
 310 of the for loop we use the shorthand  $G^{(a+1)} = G^{(a)} \setminus C_a^\#$  to denote the graph formed from

311  $G^{(a)}$  by removing the vertices  $C_a^\#$ . We do not analyze the theoretical performance of ICP+RWT  
 312 here<sup>2</sup> but we provide numerical evidence that ICP+RWT is competitive with state-of-the-art  
 313 semi-supervised graph clustering algorithms in §10.3.

---

**Algorithm 6.1** ICP+RWT
 

---

**Input:** Adjacency matrix  $A$ , labeled data  $\Gamma_a \subset C_a$  for  $a = 1, \dots, k$ . Parameters  $\epsilon \in (0, 1)$ ,  
 $R \in [0, 1)$ ,  $\hat{n}_a \approx n_a$  and  $s_a \approx \epsilon n_a$  for  $a = 1, \dots, k$ , and  $t \in \mathbb{Z}_+$

**Initialize:**  $G^{(1)} = G$  and  $A^{(1)} = A$ .

**for**  $a = 1, \dots, k$  **do**

Let  $C_a^\# = \text{CP+RWT}(A^{(a)}, \Gamma_a, \epsilon, R, s_a, \hat{n}_a, t)$

Let  $G^{(a+1)} = G^{(a)} \setminus C_a^\#$  and let  $A^{(a+1)}$  be the adjacency matrix of  $G^{(a+1)}$ .

**end for**

**Output:**  $C_1^\#, \dots, C_k^\#$

---

314 **7. Computational Complexity.** In this section we discuss the run times of the algorithms  
 315 introduced in this paper. Let  $\mathcal{T}_m$  denote the cost of a matrix-vector multiply with  $A$ ,  $L$  or  $P$   
 316 (they are all of the same magnitude).

317 **Theorem 7.1.** *RWThresh* requires  $O(n \log(n) + t\mathcal{T}_m)$  operations, where  $t$  is the depth of the  
 318 random walk.

319 *Proof.* Computing  $\mathbf{v}^{(t)}$  requires  $t$  matrix-vector multiplies and hence requires  $O(t\mathcal{T}_m)$  op-  
 320 erations. Sorting  $\mathbf{v}^{(t)}$  in order to find  $\Omega$  requires  $O(n \log(n))$  operations. ■

321 Let us now analyze the complexity of **ClusterPursuit**

322 **Theorem 7.2.** *ClusterPursuit* requires  $O(\mathcal{T}_m \log(n))$  operations.

323 *Remark 7.3.* Note that if  $A$  is stored as a sparse matrix then  $\mathcal{T}_m = O(nd_{\max})$  in which  
 324 case the run time of **ClusterPursuit** becomes  $O(nd_{\max} \log(n))$ .

*Proof.* The run time of **ClusterPursuit** is dominated by the cost of the call to **SubspacePursuit**■  
 (see Algorithm 2.1) in step (3) which costs  $m$  times the cost of each iteration. We now bound  
 the cost of each iteration. The cost of the  $j$ -th iteration is dominated by the cost of solving  
 the least squares problem:

$$\arg \min_{\mathbf{z} \in \mathbb{R}^n} \left\{ \|L\mathbf{z} - \mathbf{y}\|_2 : \text{supp}(\mathbf{z}) \subset \hat{S}^j \right\}.$$

325 (step (4) in the “for” loop of Algorithm 2.1). Because of the support condition, and because  
 326  $|\hat{S}^j| = 2s \leq 0.26n_1$ , this is equivalent to the least squares problem:

327 (7.1) 
$$\arg \min_{\mathbf{z} \in \mathbb{R}^{2s}} \left\{ \|L_{\hat{S}^j} \mathbf{z} - \mathbf{y}\|_2 \right\}$$

328 We recommend using an iterative method, such as conjugate gradient (in our implementation  
 329 we use MATLAB’s `lsqr` operation). Fortunately, as pointed out in [39], the matrix in question,

---

<sup>2</sup>There is a minor technical difficulty: one needs to show that if  $G$  is drawn from a model satisfying assumptions (A1)–(A4) then each  $G^{(a)}$  is also drawn from a model satisfying assumptions (A1)–(A4).

330  $L_{\hat{S}_j}$  is extremely well conditioned. This is because  $\delta_{2s}(L) \leq \delta_{3s}(L) \leq 0.45$ , as shown in the  
 331 proof of Theorem 3.4. By [39], specifically Proposition 3.1 and the discussion of §5, this  
 332 implies that the condition number is small:

$$333 \quad \kappa(L_{\hat{S}_j}^\top L_{\hat{S}_j}) := \frac{\lambda_{\max}(L_{\hat{S}_j}^\top L_{\hat{S}_j})}{\lambda_{\min}(L_{\hat{S}_j}^\top L_{\hat{S}_j})} \leq \frac{1 + \delta_{2s}}{1 - \delta_{2s}} \leq 2.64$$

334 The upshot of this is that it only requires a constant number of iterations of conjugate gra-  
 335 dient to approximate the solution to the least-squares Problem (7.1) to within an acceptable  
 336 tolerance. Indeed, Corollary 5.3 of [39] argues that three iterations suffices. We play it safe  
 337 by performing ten iterations. The cost of each iteration of conjugate gradient is equal to (a  
 338 constant times) the cost of a matrix vector multiply by  $L_{\hat{S}_j}$  or  $L_{\hat{S}_j}^\top$ , which is  $\mathcal{T}_m$ . Hence the  
 339 total cost of step (3) of `ClusterPursuit` is  $O(m\mathcal{T}_m) = O(\log(n)\mathcal{T}_m)$  because we are taking  
 340  $m = O(\log(n))$ . ■

341 As a direct consequence of Theorems 7.1 and 7.2, we get that `CP+RWT` runs in time  
 342  $O((nd_{\max} \log(n)))$ . If the number of clusters,  $k$ , is  $O(1)$ , we get that `ICP+RWT` also runs in  
 343 time  $O((nd_{\max} \log(n)))$ .

344 **8. Comparison with Existing Literature.** `ClusterPursuit` can naturally be compared  
 345 with other cut improvement algorithms such as `FlowImprove` [2], `LocalFlow` [41] and `SimpleLocal` [50].  
 346 We note that the performance guarantees for these three algorithms are of a different  
 347 flavor to ours. Specifically, and translating into the notation of this paper, they bound the  
 348 conductance of the improved cut,  $C_1^\#$ , by some function of the original cut,  $\Omega$ . In contrast,  
 349 our performance guarantees for `ClusterPursuit` are of a more statistical nature. In terms  
 350 of run-time, `LocalFlow` and `SimpleLocal` are strongly local, so have run times  $O(\text{vol}(\Omega)^\alpha)$   
 351 for  $\alpha \geq 1$ . While this is certainly better than `ClusterPursuit` for finding small clusters, *ie*  
 352 when  $|\Omega| = O(1)$ , these run times become less attractive for even moderate sized clusters, *eg*  
 353  $|C_1| = O(\sqrt{n})$ . In §10 we demonstrate that `ClusterPursuit` is several orders of magnitude  
 354 faster than `FlowImprove` and `SimpleLocal` in the regime  $|C_1| = O(n)$ .

355  
 356 The idea of combining a fast, diffusion based clustering algorithm with a refinement pro-  
 357 cedure to create a local clustering algorithm is not new. See, for example, the algorithms  
 358 `LEMON` [27, 35], `LOSP` and `LOSP++` [34], `LBSA` [44], and `FlowSeed` [51]. We compare `CP+RWT` to a  
 359 selection of these algorithm in §10. We note that there exist many diffusion-based local clus-  
 360 tering algorithms that may find better approximations to  $C_1$  than `RWThresh`. See for example,  
 361 `PPR-Grow` [6], `HK-Grow` [31] or `CapacityReleasingDiffusion` [52]. We emphasize that the  
 362 main advantage of `RWThresh` is that it rapidly and provably finds good enough initial cuts,  $\Omega$ ,  
 363 to be fed into `ClusterPursuit`. We show in §10 that the combination `CP+RWT` typically out-  
 364 performs these diffusion-only approaches, particularly for large, sparsely connected clusters.

365  
 366 The analysis of `CP+RWT` contained in §3–5 can be compared to the recent works [52] and  
 367 [26]. In both the performance of a local clustering algorithm on graphs drawn from a certain  
 368 probabilistic model is studied. In both papers, the model is more general in one sense: there  
 369 is no restriction on the structure of  $V \setminus C_1$ , but more restrictive in other senses: the ratio

370  $d^{\text{out}}/d^{\text{in}}$  must be at most  $O(1/\log^2(n_1))$  in [52] while the results in [26] are most meaningful  
 371 when  $n_1 = O(1)$  and  $d^{\text{out}} = O(1)$ . In contrast, our results tackle the regime where  $n_1 = O(n)$   
 372 and  $d^{\text{out}}/d^{\text{in}}$  can be bounded by an arbitrarily slowly decaying function of  $n_1$ .

373

374 Finally, we mention several recent works that combine notions of sparsity and local clus-  
 375 tering. In particular, we mention the works of Fountoulakis, Gleich, Mahoney *et al* [25, 22, 26]  
 376 which introduce and study the  $\ell_1$  regularized page rank problem. The algorithms LOSP and  
 377 LOSP++ also set up and solve a sparse recovery problem, although with an additional non-  
 378 negativity requirement. However, to the best of the authors' knowledge, ClusterPursuit is  
 379 the first algorithm that explicitly phrases the problem of improving a cut,  $\Omega$ , as the problem  
 380 of finding a sparse change to the indicator vector  $\mathbf{1}_\Omega$ .

381 **9. Which Probabilistic Models Satisfy our Assumptions?** First, we verify that a well-  
 382 studied model of graphs with clusters, namely the stochastic block model, satisfies Assump-  
 383 tions (A1)–(A4) of §2.3. We first remind the reader of the simpler Erdős - R enyi model:

384 **Definition 9.1.** We say  $G = (V, E)$  is drawn from the Erdős - R enyi model on  $n$  vertices  
 385 with parameter  $p$  (and write  $G \sim ER(n, p)$ ) if  $V = [n]$  and  $\mathbb{P}[\{i, j\} \in E] = p$  for  $i, j \in V$ , with  
 386 all such probabilities being independent.

387 **Definition 9.2** ([28, 3]). Let  $\mathbf{n} = (n_1, \dots, n_k)$  be a vector of positive integers, and let  $P$  be  
 388 a  $k \times k$  symmetric matrix with entries  $P_{ab} \in [0, 1]$  for all  $a, b$ . We say a graph  $G = (V, E)$   
 389 is drawn from the Stochastic Block Model (written  $G \sim SBM(\mathbf{n}, P)$ ) if there exists a partition  
 390  $V = C_1 \cup C_2 \dots \cup C_k$  with  $|C_a| = n_a$  such that any vertices  $i \in C_a$  and  $j \in C_b$  are connected  
 391 by an edge with probability  $P_{ab}$ , and all edges are inserted independently.

392 Note that if  $G \sim SBM(\mathbf{n}, P)$  then each  $G_{C_a} \sim ER(n_a, P_{aa})$ . Without loss of generality,  
 393 we shall assume that  $n_1 \leq n_2 \leq \dots \leq n_k$ . In an appendix, we shall prove the following:

394 **Theorem 9.3.** Suppose that  $n_1 = O(n) \rightarrow \infty$ ,  $P_{aa} = \omega \log(n)/n_a$  for any  $\omega \rightarrow \infty$  and  
 395  $P_{ab} = (\beta + o(1)) \log(n)/n$  for any  $a \neq b$  where  $\beta \geq 0$  is a constant. Then  $SBM(\mathbf{n}, P)$  satisfies  
 396 assumptions (A1)–(A4).

397 *Proof.* See Appendix C. ■

398 As a consequence of this theorem we have that, given a small fraction of vertices in  $C_1$ ,  
 399 CP+RWT will reliably return a  $C_1^\#$  with  $|C_1 \triangle C_1^\#| = o(n_1)$ . We experimentally confirm this  
 400 in §10 for  $\omega \sim \log(n)$ . In this regime we have that  $d_{\text{max}} = O(\log^2(n))$  with high probability,  
 401 hence the run time of CP+RWT is  $O(n \log^3(n))$  by Theorem 7.2.

402

403 It is interesting to contrast this result with what is known for the global clustering problem  
 404 for the stochastic block model. There are several unsupervised algorithms, see for example [4]  
 405 and [38], that return a partition  $V = C_1^\# \cup C_2^\# \cup \dots \cup C_k^\#$  such that  $C_a^\# = C_a$  with high  
 406 probability. However these approaches either have impractically high run times [38] or are  
 407 tricky to implement in practice [4]. In contrast, CP+RWT has a low run time, in theory and  
 408 in practice, and can be implemented in a few lines of code. In addition, the ‘‘one cluster  
 409 at a time’’ nature of CP+RWT affords an additional flexibility that may be useful in certain  
 410 circumstances.

411

412 On the other hand, we have had less success with using CP+RWT for certain random geo-  
 413 metric graphs arising as  $K$ -NN graphs of point clouds in  $\mathbb{R}^d$ . We note that CP+RWT is most  
 414 effective when the adjacency matrix of the  $K$ -NN graph is sparse but has its non-zero en-  
 415 tries uniformly distributed. In contrast, for certain artificial data sets, for example points  
 416 drawn from a thickened line or sphere embedded in a high dimensional space, this adjacency  
 417 matrix tends to exhibit a banded structure—at least when nearest neighbors are determined  
 418 using the Euclidean metric. Experimentally, we have observed that CP+RWT performs poorly  
 419 on these data sets. However, this problem is to a large extent particular to the use of the  
 420 Euclidean metric. In particular, when a data-driven metric such as those detailed in [36] is  
 421 used to construct the  $K$ -NN graph, CP+RWT performs much better. Moreover, even when using  
 422 the Euclidean metric CP+RWT still performs extremely well on real data sets, such as MNIST,  
 423 COIL and Optdigits, which are frequently thought of as consisting of data points drawn from  
 424 a low-dimensional manifold embedded in a high dimensional space (see §10.3).

425

**10. Numerical Experiments.** We compare the algorithms ClusterPursuit, CP+RWT and  
 426 ICP+RWT to the state of the art on the various problems they are designed to solve. Specifically,  
 427 in §10.1 we compare the performance of ClusterPursuit on the cut improvement task to  
 428 two baseline algorithms, namely FlowImprove and SimpleLocal, for graphs drawn from the  
 429 stochastic block model. We also compare CP+RWT to the local clustering algorithms HK-Grow,  
 430 PPR-Grow and LBSA for the same data.<sup>3</sup> In §10.2 we repeat this experiment for social networks.  
 431 We take care to choose our data sets and performance measures to allow for easy comparison  
 432 with similar work in [52]. In §10.3 we test the performance of ICP+RWT on two data sets  
 433 commonly studied in the machine learning community—MNIST and OptDigits. We provide  
 434 a detailed description of the implementation of all algorithms considered in the supplementary  
 435 material.

436

**10.1. Synthetic Data Sets.** We consider graphs drawn from  $\text{SBM}(\mathbf{n}^{(i)}, P^{(i)})$  for two  
 437 different sets of parameters. The first set:  $\mathbf{n}^{(1)} = (n_1, 1.5n_1, 2.5n_1, 5n_1)$  and  $P^{(1)}$  with  
 438  $P_{aa} = \log^2(n)/2$  and  $P_{ab} = 5\log(n)/n$  for all  $a \neq b$  is designed to satisfy the conditions  
 439 of Theorem 9.3 while presenting a challenge to existing clustering algorithms. The second set:  
 440  $\mathbf{n}^{(2)} = (n_1, 10n_1)$  and  $P^{(2)} = \begin{bmatrix} 2\log^2(n)/n & \log(n)/n \\ \log(n)/n & \log(n)/n \end{bmatrix}$  goes beyond the assumptions of Theorem  
 441 9.3 and is essentially the planted cluster model studied in [26] and elsewhere. For both sets  
 442 of parameters we perform two experiments. In the first we test the performance of the three  
 443 cut improvement algorithms when initialized with an  $\Omega$  “close” to  $C_1$ . This  $\Omega$  is found using  
 444 RWThresh. In the second we compare the performance of CP+RWT with the performance of the  
 445 local clustering algorithms mentioned above. For both experiments we report both run time  
 446 and accuracy, as measured by the Jaccard Index in Figure 2 and in Figure 3, respectively.

447

**10.2. Social Networks.** The well-known facebook100 dataset consists of anonymized  
 448 Facebook friendship networks at 100 American universities, and was first introduced and  
 449 studied in [48]. Certain demographic markers (year of entry, residence etc.) were also collected

---

<sup>3</sup>While there are certainly other worthy local clustering algorithms that deserve to be included, such as CapacityReleasingDiffusion [52] and FlowSeed [51], we stick to algorithms with a freely available MATLAB implementation

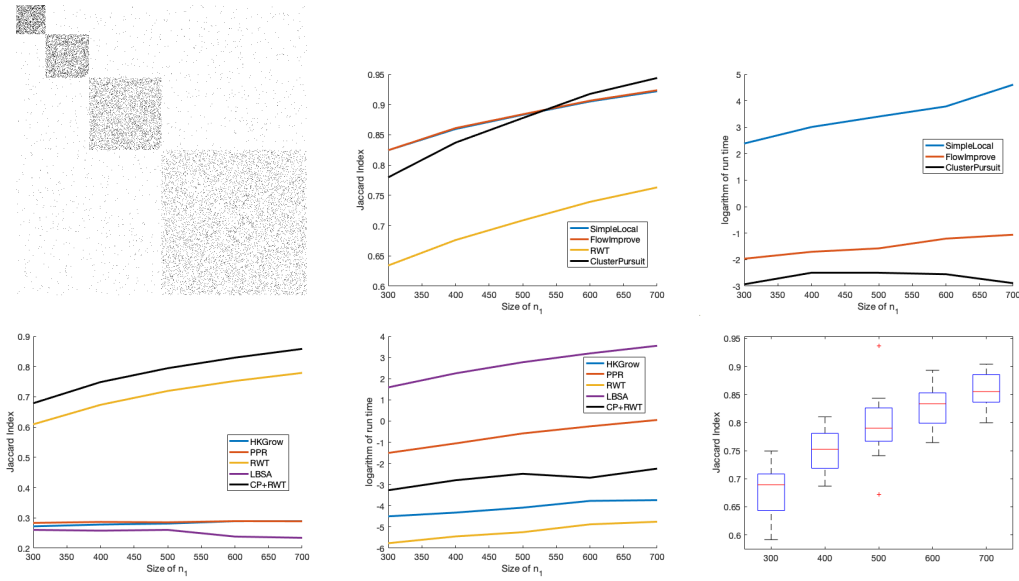


Figure 2: *Top row, left to right:* Stylized representation of the adjacency matrix of graphs drawn from  $\text{SBM}(\mathbf{n}^{(1)}, P^{(1)})$ , Jaccard index for results of cut improvement (**SimpleLocal** and **FlowImprove** always have the same Jaccard index) and (log. of) run time for the three cut improvement algorithms. Note that **ClusterPursuit** is at least an order of magnitude faster than the other two, even though **FlowImprove** is implemented in *C*. *Bottom row, left to right:* Jaccard index for local clustering (The poor performance of the other methods is not an implementation issue. Rather, it is a consequence of the small gap between  $P_{aa}^{(1)}$  and  $P_{ab}^{(1)}$ ). (Log. of) run time for local clustering. Box plot of Jaccard index for **CP+RWT**.

450 in an anonymized format. One can think of vertices sharing the same marker as defining a  
 451 ground truth cluster, although some of these clusters are extremely noisy. We focus on four  
 452 clusters identified in [52] as having good (*ie* low) or moderately good conductance scores,  
 453 namely Johns Hopkins class of 2009, Rice University dorm 203, Simmons College class of  
 454 2009 and Colgate University class of 2006. The details of these clusters are displayed in Table  
 455 1. For ease of comparison with the results of [52] we report accuracy using precision and recall  
 456 scores. We remind the reader that, in the notation of this paper, precision =  $|C_1 \cap C_1^\#|/|C_1^\#|$   
 457 and recall =  $|C_1 \cap C_1^\#|/|C_1|$ . It is desirable to have both of these values as close to 1 as  
 458 possible. For all four experiments we take  $\Gamma$  to be selected uniformly and at random from  
 459  $C_1$ , with  $|\Gamma| = 0.02n_1$ . We average over fifty independent trials. There results are shown in  
 460 Figure 4.

461 **10.3. Machine Learning Benchmarks.** We consider two venerable benchmark data sets:  
 462

463 **OptDigits.** This data set consists of grayscale images of handwritten digits 0–9 of size  
 464  $8 \times 8$ . There are  $n = 5620$  images and the clusters are fairly well balanced with approximately  
 465 560 images of each digit.



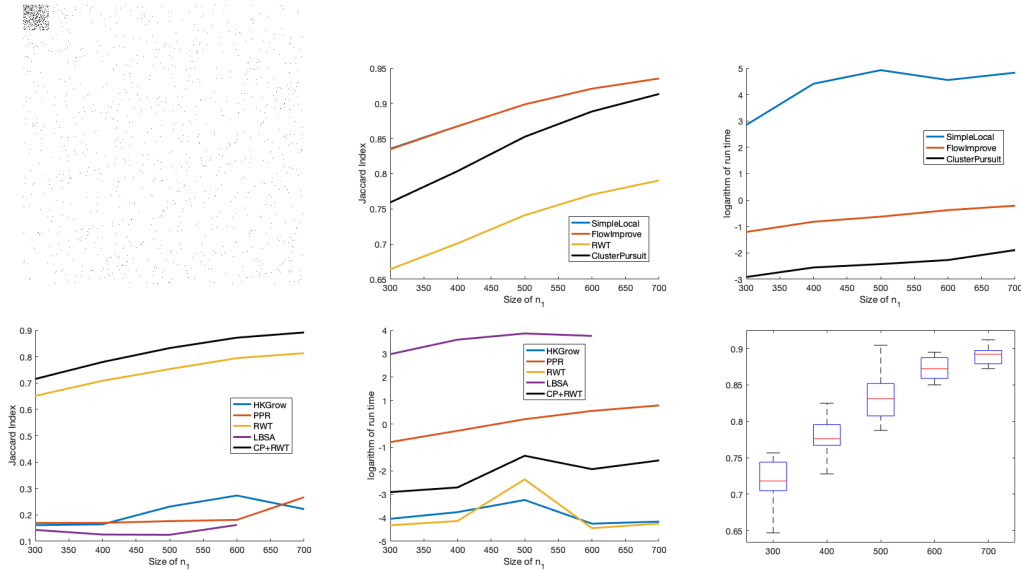


Figure 3: *Top row, left to right:* Stylized representation of the adjacency matrix of graphs drawn from  $\text{SBM}(\mathbf{n}^{(2)}, P^{(2)})$ , Jaccard index for results of cut improvement (again, **SimpleLocal** and **FlowImprove** always have the same Jaccard index). (log. of) Run time for the three cut improvement algorithms. *Bottom row, left to right:* Jaccard index for local clustering (Again, the poor performance of the benchmark methods is a consequence of the challenging SBM parameters chosen). (Log. of) run time for local clustering. Box plot of Jaccard index for **CP+RWT**.

School	Cluster	Size of graph	Size of Cluster	Conductance
Johns Hopkins	Class of 2009	5180	910	0.21
Rice	Dorm. 203	4087	406	0.47
Simmons	Class of 2009	1518	289	0.11
Colgate	Class of 2006	3482	557	0.49

Table 1: Basic properties of the four social networks studied.

466 **MNIST.** This data set also consists of grayscale images of the handwritten digits 0–9  
 467 although here there are  $n = 70\,000$  images, all of size  $28 \times 28$ . There are approximately 7 000  
 468 images of each digit.

469

470 For each data set we form a  $k$ -NN graph using the procedure presented in [30] and described  
 471 in detail in the supplementary material. The labeled data,  $\Gamma_a$ , was sampled uniformly at  
 472 random from  $C_a$ , and each is of size  $g|C_a|$ . The accuracy of the classification given by **ICP+RWT**,  
 473 for increasing  $g$ , is presented in Table 2. All results are averaged over twenty independent  
 474 trials.

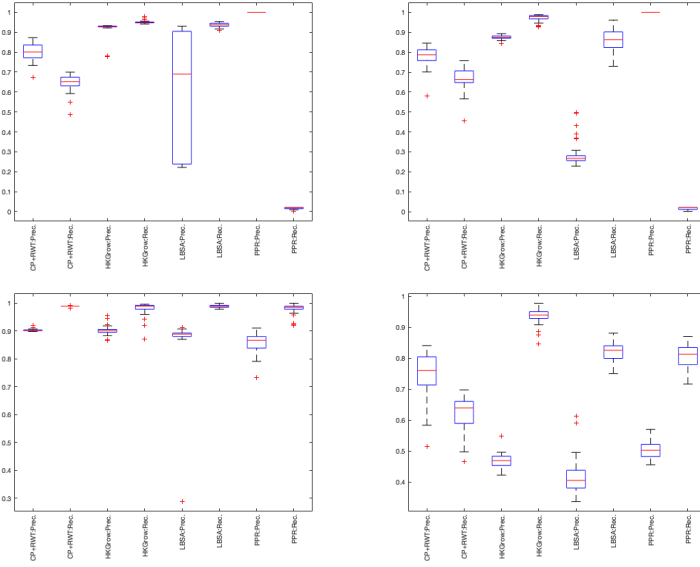


Figure 4: Precision and Recall for various local clustering algorithms on the social networks described in Table 1. Clockwise from top left: Johns Hopkins, Rice, Colgate and Simmons. Note that CP+RWT consistently achieves high precision without sacrificing recall.

Data Set	% Labeled Data				
	0.5	1	1.5	2	2.5
MNIST	96.41%	97.32%	97.44%	97.52%	97.50%
OptDigits	91.88%	95.47%	97.16%	98.06%	98.08%

Table 2: Classification accuracy, as a function of amount of labeled data, for ICP+RWT on two well-studied benchmark data sets.

Method	Labeled	Accuracy
TVRF [54]	600	96.8%
ICP+RWT	700	97.32%
Multi-Class MBO with Auction Dynamics [30]	700	97.43%
ICP+RWT	1050	97.44%
Ladder Networks [42]	1000	99.16%

Table 3: Comparing ICP+RWT to other, state-of-the-art, semi-supervised methods on MNIST. TVRF and Multi-Class MBO are graph-based, and have similar run times to ICP+RWT. The Ladder Network approach uses a deep neural network and hence requires training ( $\sim 2$  hours on a GPU) before it can be used for classification.

475

## REFERENCES

476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527

- [1] Reid Andersen, Kevin J. Lang, Communities from seed sets, Proceedings of the 15th International Conference on World Wide Web, 223–232, 2006.
- [2] Reid Andersen, Kevin J. Lang, An algorithm for improving graph partitions, Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms, 2008.
- [3] Emmanuel Abbe. Community detection and stochastic block models: Recent developments. The Journal of Machine Learning Research, 18(1):6446–6531, 2017.
- [4] Emmanuel Abbe and Colin Sandon. Recovering communities in the general stochastic block model without knowing the parameters. In Advances in Neural Information Processing Systems, 676–684, 2015.
- [5] 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, 36–43, 2005.
- [6] Reid Andersen, Fan Chung, and Kevin Lang. Using pagerank to locally partition a graph. Internet Mathematics, 4(1):35–64, 2007.
- [7] Béla Bollobás. Vertices of given degree in a random graph. Journal of Graph Theory, 6(2):147–155, 1982.
- [8] Béla Bollobás. Random graphs, Cambridge University Press, 2001.
- [9] Emmanuel J. Candes, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. IEEE Transactions on Information Theory, 52(2):489–509, 2006.
- [10] Olivier Chappelle, Bernard Schölkopf and Alexander Zien. Semi-Supervised Learning. MIT Press. 2006.
- [11] Fan Chung. Spectral graph theory (CBMS Regional Conference Series in Mathematics, no. 92). 1996.
- [12] Fan Chung. The heat kernel as the pagerank of a graph. Proceedings of the National Academy of Sciences, 104(50):19735–19740, 2007.
- [13] Fan Chung. Random walks and local cuts in graphs. Linear Algebra and its Applications, 423(1):22–32, 2007.
- [14] Fan Chung and Ron Graham. Quasi-random graphs with given degree sequences. Random Structures & Algorithms, 32(1):1–19, 2008.
- [15] Fan Chung and Mary Radcliffe, On the spectra of general random graphs. The Electronic Journal of Combinatorics, 18(1):215, 2011.
- [16] Fan Chung and Olivia Simpson. Computing heat kernel pagerank and a local clustering algorithm. European Journal of Combinatorics, 68:96–119, 2018.
- [17] Wei Dai and Olgica Milenkovic. Subspace pursuit for compressive sensing signal reconstruction. IEEE Transactions on Information Theory, 55(5):2230–2249, 2009.
- [18] Pasquale De Meo, Emilio Ferrara, Giacomo Fiumara, and Alessandro Provetti. Generalized Louvain method for community detection in large networks. In 2011 11th International Conference on Intelligent Systems Design and Applications: 88–93, 2011.
- [19] Inderjit S. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning. In Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining:269–274, 2001.
- [20] David L. Donoho. Compressed sensing. IEEE Transactions on Information Theory, 52(4):1289–1306, 2006.
- [21] Simon Foucart and Holger Rauhut. A Mathematical Introduction to Compressive Sensing. Springer Science & Business Media, 2013.
- [22] Kimon Fountoulakis, David Gleich and Michael Mahoney. An optimization approach to locally-biased graph algorithms. Proceedings of the IEEE 105.2: 256-272, 2017.
- [23] Alan Frieze and Michal Karoński. Introduction to random graphs. Cambridge University Press, 2016.
- [24] Michelle Girvan and Mark E.J. Newman. Community structure in social and biological networks. Proceedings of the National Academy of Sciences, 99(12):7821–7826, 2002.

- 528 [25] David Gleich and Michael Mahoney. Anti-differentiating approximation algorithms: A case  
529 study with min-cuts, spectral and flow. International Conference on Machine Learning,  
530 2014.
- 531 [26] Wooseok Ha, Kimon Fountoulakis, and Michael W. Mahoney. Statistical guarantees for local  
532 graph clustering. arXiv preprint arXiv:1906.04863 (2019).
- 533 [27] Kun He, Yiwei Sun, David Bindel, John Hopcroft, and Yixuan Li. Detecting overlapping  
534 communities from local spectral subspaces. In 2015 IEEE International Conference on  
535 Data Mining: 769–774, 2015.
- 536 [28] Paul Holland, Kathryn Blackmond Laskey and Samuel Leinhardt. Stochastic blockmodels:  
537 First steps. Social Networks, 5(2): 109–137, 1983.
- 538 [29] Matthew A. Herman and Thomas Strohmer. General deviants: An analysis of perturbations in  
539 compressed sensing. IEEE Journal of Selected Topics in Signal Processing, 4(2):342–349,  
540 2010.
- 541 [30] Matt Jacobs, Ekaterina Merkurjev, and Selim Esedođlu. Auction dynamics: A volume con-  
542 strained MBO scheme. Journal of Computational Physics, 354:288–310, 2018.
- 543 [31] Kyle Kloster and David F. Gleich. Heat kernel based community detection. In Proceedings  
544 of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data  
545 Mining: 1386–1395, 2014.
- 546 [32] Jeffrey Lewis, Keith Poole, Howard Rosenthal, Adam Boche, Aaron Rudkin, and Luke Sonnet.  
547 Voteview: Congressional Roll-Call Votes Database. <https://voteview.com/>
- 548 [33] Haifeng Li. Improved analysis of SP and CoSaMP under total perturbations. EURASIP Journal  
549 on Advances in Signal Processing, 2016(1):112, 2016.
- 550 [34] Yixuan Li, Kun He, David Bindel, and John E. Hopcroft. Uncovering the small community  
551 structure in large networks: A local spectral approach. In Proceedings of the 24th Interna-  
552 tional Conference on World Wide Web: 658–668, 2015.
- 553 [35] Yixuan Li, Kun He, Kyle Kloster, David Bindel, and John E. Hopcroft. Local spectral clustering  
554 for overlapping community detection. ACM Transactions on Knowledge Discovery from  
555 Data (TKDD), 12(2):17, 2018.
- 556 [36] Daniel McKenzie and Steven Damelin. Power weighted shortest paths for clustering Euclidean  
557 data. Foundations of Data Science, 1.3 (2019).
- 558 [37] Michael W. Mahoney, Lorenzo Orecchia, and Nisheeth K. Vishnoi. A local spectral method for  
559 graphs: With applications to improving graph partitions and exploring data graphs locally.  
560 Journal of Machine Learning Research, 13(Aug):2339–2365, 2012.
- 561 [38] Elchanan Mossel, Joe Neeman and Allan Sly. A proof of the block model threshold conjecture.  
562 Combinatorica 38.3: 665–708, 2018.
- 563 [39] Deanna Needell and Joel A. Tropp. CoSaMP: Iterative signal recovery from incomplete and  
564 inaccurate samples. Applied and Computational Harmonic Analysis, 26(3):301–321, 2009.
- 565 [40] Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an  
566 algorithm. In Advances in Neural Information Processing Systems: 849–856, 2002.
- 567 [41] Lorenzo Orecchia, Zeyuan Allen Zhu. Flow-based algorithms for local graph clustering. Pro-  
568 ceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms: 1267–  
569 1286, 2014.
- 570 [42] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semisu-  
571 pervised learning with ladder networks. In Advances in Neural Information Processing  
572 Systems: 3546–3554, 2015.
- 573 [43] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. IEEE Transactions  
574 on Pattern Analysis and Machine Intelligence, 22(8):888–905, 2000.
- 575 [44] Pan Shi, Kun He, David Bindel, and John E. Hopcroft. Locally-biased spectral approximation  
576 for community detection. Knowledge-Based Systems, 164:459–472, 2019.
- 577 [45] Daniel A. Spielman and Shang-Hua Teng. Nearly linear time algorithms for graph partitioning,  
578 graph sparsification, and solving linear systems, In Proceedings of the STOC (4), 2004.
- 579 [46] Daniel A. Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs  
580 and its application to nearly linear time graph partitioning. SIAM Journal on Computing,  
581 42(1):1–26, 2013.

- 582 [47] Robert C. Thompson. Principal submatrices IX: Interlacing inequalities for singular values of  
583 submatrices. *Linear Algebra and its Applications*, 5(1):1–12, 1972.
- 584 [48] Amanda L. Traud, Peter J. Mucha, and Mason A. Porter. Social structure of Facebook net-  
585 works. *Physica A: Statistical Mechanics and its Applications*, 391(16):4165–4180, 2012.
- 586 [49] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–  
587 416, 2007.
- 588 [50] Nate Veldt, David F. Gleich, Micheal W. Mahoney, A simple and strongly-local flow-based  
589 method for cut improvement, *International Conference on Machine Learning* (2016).
- 590 [51] Nate Veldt, Christine Klymko, David F. Gleich, Flow-based local graph clustering with better  
591 seed set inclusion, *Proceedings of the 2019 SIAM International Conference on Data Mining*,  
592 378–386 (2019).
- 593 [52] Di Wang, Kimon Fountoulakis, Monika Henzinger, Michael W. Mahoney, Satish Rao, Capacity  
594 releasing diffusion for speed and locality, *Proceedings of the 34th International Conference  
595 on Machine Learning*, 70:3598–3607, 2017.
- 596 [53] James D. Wilson, Simi Wang, Peter J. Mucha, Shankar Bhamidi, Andrew B. Nobel, A testing  
597 based extraction algorithm for identifying significant communities in networks, *The Annals  
598 of Applied Statistics*, 8(3):1853–1891, 2014.
- 599 [54] Ke Yin and Xue-Cheng Tai. An effective region force for some variational models for learning  
600 and clustering, *Journal of Scientific Computing*, 74 (2018), 175–196.
- 601 [55] Lih Zelnik-Manor and Pietro Perona. Self-tuning spectral clustering. In *Advances in Neural  
602 Information Processing Systems*, pages 1601–1608, 2005.

603 **A. Restricted Isometry Property for Laplacians.** In this section, we prove parts (1) and  
604 (2) of Theorem 3.3. We proceed via a series of lemmas.

605 **A.1. Restricted Isometry Property for  $L^{\text{in}}$ .**

606 **Lemma A.1.** *Let  $G$  be any connected graph on  $n_0$  vertices, and let  $s < n_0$ . Let  $\lambda_i := \lambda_i(L)$   
607 denote the  $i$ -th smallest eigenvalue of  $L$ . Then:*

608 
$$\delta_s(L) \leq \max\left\{1 - \lambda_2^2 \left(\frac{d_{\min}}{d_{\max}} - \frac{d_{\max}}{d_{\min}} \frac{s}{n_0}\right), \lambda_{\max}^2 - 1\right\}.$$

609 *Proof.* Recall that the  $s$ -Restricted Isometry Constant  $\delta_s(L)$  is the smallest  $\delta$  such that,  
610 for any  $\mathbf{v}$  with  $\|\mathbf{v}\|_0 \leq s$  and  $\|\mathbf{v}\|_2 = 1$ :  $(1 - \delta) \leq \|L\mathbf{v}\|_2^2 \leq (1 + \delta)$ . The RHS bound is  
611 straightforward since

612 
$$\|L\mathbf{v}\|_2 \leq \|L\|_2 \|\mathbf{v}\|_2 = \lambda_{\max} \mathbf{1} = \lambda_{\max}.$$

613 The LHS bound requires some work. Recall that  $L = I - D^{-1}A$ . This matrix is not symmetric,  
614 but  $L^{\text{sym}} = I - D^{-1/2}AD^{-1/2}$  is. By Lemma 2.3  $L$  and  $L^{\text{sym}}$  have the same eigenvalues. Let  
615  $\mathbf{w}_1, \dots, \mathbf{w}_{n_0}$  be an orthonormal eigenbasis for  $L^{\text{sym}}$ . These eigenvectors are well studied (see,  
616 for example, [11]) and in particular  $\mathbf{w}_1 = \frac{1}{\sqrt{\text{vol}(G)}}D^{1/2}\mathbf{1}$  where  $\mathbf{1}$  is the all-ones vector. Observe  
617 that:

618 
$$L\mathbf{v} = D^{-1/2} \left( D^{1/2}LD^{-1/2} \right) D^{1/2}\mathbf{v} = D^{-1/2}L^{\text{sym}}D^{1/2}\mathbf{v} = D^{-1/2}L^{\text{sym}}\mathbf{z},$$

619 where  $\mathbf{z} := D^{1/2}\mathbf{v}$ . It follows that:

620 (A.1) 
$$\|L\mathbf{v}\|_2 = \|D^{-1/2}L^{\text{sym}}\mathbf{z}\|_2 \geq \frac{1}{\sqrt{d_{\max}}} \|L^{\text{sym}}\mathbf{z}\|_2.$$

621 Express  $\mathbf{z}$  in terms of the orthonormal basis  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ , namely  $\mathbf{z} = \sum_{i=1}^{n_0} \alpha_i \mathbf{w}_i$ . Then:

$$622 \quad \|L^{\text{sym}} \mathbf{z}\|_2^2 = \left\| \sum_{i=1}^{n_0} \alpha_i \lambda_i \mathbf{w}_i \right\|_2^2 = \left\| \sum_{i=2}^{n_0} \alpha_i \lambda_i \mathbf{w}_i \right\|_2^2 \geq \lambda_2^2 \left( \sum_{i=2}^{n_0} \alpha_i^2 \right)$$

623 and  $\sum_{i=2}^{n_0} \alpha_i^2 = \|\mathbf{z}\|_2^2 - \alpha_1^2$ . We now bound  $\|\mathbf{z}\|_2$  and  $\alpha_1$ .

$$624 \quad \|\mathbf{z}\|_2^2 = \|D^{1/2} \mathbf{v}\|_2^2 \geq \left( \sqrt{d_{\min}} \right)^2 \|\mathbf{v}\|_2^2 = d_{\min}$$

625 while:

$$626 \quad \alpha_1 = \langle \mathbf{z}, \mathbf{w}_1 \rangle = \langle D^{1/2} \mathbf{v}, \frac{1}{\sqrt{\text{vol}(G)}} D^{1/2} \mathbf{1} \rangle = \frac{1}{\sqrt{\text{vol}(G)}} \langle \mathbf{v}, D \mathbf{1} \rangle \leq \frac{d_{\max}}{\sqrt{\text{vol}(G)}} \langle \mathbf{v}, \mathbf{1} \rangle.$$

627 We now use the assumptions on  $\mathbf{v}$ . Specifically  $\langle \mathbf{v}, \mathbf{1} \rangle \leq \|\mathbf{v}\|_1 \leq \sqrt{s} \|\mathbf{v}\|_2 = \sqrt{s}$  and so

$$628 \quad \alpha_1 \leq d_{\max} \frac{\sqrt{s}}{\sqrt{\text{vol}(G)}} \leq d_{\max} \frac{\sqrt{s}}{\sqrt{d_{\min} n_0}} = \frac{d_{\max}}{\sqrt{d_{\min}}} \frac{\sqrt{s}}{\sqrt{n_0}}.$$

629 Returning to equation (A.1):

$$630 \quad \|L \mathbf{v}\|_2^2 \geq \frac{1}{d_{\max}} \|L^{\text{sym}} \mathbf{z}\|_2^2 \geq \frac{1}{d_{\max}} \lambda_2^2 \left( d_{\min} - \frac{d_{\max}^2 s}{d_{\min} n_0} \right) = \lambda_2^2 \left( \frac{d_{\min}}{d_{\max}} - \frac{d_{\max} s}{d_{\min} n_0} \right).$$

631 These yield the desired estimate. ■

632 **Theorem A.2.** Let  $G \sim \mathcal{G}_n$  with  $\mathcal{G}_n$  satisfying (A2) and (A4). Then for any  $\gamma \in (0, 1)$ , we

633 have that  $\delta_{\gamma n_a}(L^{\text{in}}) \leq \frac{n_a}{n_1} \gamma + o(1)$ .

634 *Proof.* Firstly, observe that  $L^{\text{in}}$  is block diagonal with blocks  $L_{G_{C_b}}$ . For any block diagonal

635 matrix we have that  $\delta_s(L^{\text{in}}) = \max_b \delta_s(L_{G_{C_b}})$ . By Lemma A.1 we have that:

$$636 \quad (A.2) \quad \delta_s(L_{G_{C_b}}) \leq \max_b \left\{ 1 - \lambda_2(L_{G_{C_b}})^2 \left( \frac{d_{\min}^{\text{in}}}{d_{\max}^{\text{in}}} - \frac{d_{\max}^{\text{in}} s}{d_{\min}^{\text{in}} n_b} \right), \lambda_{\max}(L_{G_{C_b}})^2 - 1 \right\}.$$

637 From assumption (A4) we get that:

$$638 \quad \frac{d_{\min}^{\text{in}}}{d_{\max}^{\text{in}}} = \frac{1 - \epsilon_3}{1 + \epsilon_3} = 1 - o(1) \quad \text{and} \quad \frac{d_{\max}^{\text{in}}}{d_{\min}^{\text{in}}} = \frac{1 + \epsilon_3}{1 - \epsilon_3} = 1 + o(1).$$

639 From assumption (A2) we get that:

$$640 \quad \lambda_2(L_{G_{C_b}})^2 \geq (1 - \epsilon_1)^2 = 1 - 2\epsilon_1 + \epsilon_1^2 = 1 - o(1)$$

641 and similarly  $\lambda_{\max}(L_{G_{C_b}})^2 - 1 = o(1)$ . Plugging this in to (A.2) with  $s = \gamma n_a$  gives:

$$642 \quad \delta_{\gamma n_a}(L_{G_{C_b}}) \leq \max \left\{ \frac{\gamma n_a}{n_b} + o(1), o(1) \right\} \leq \gamma \frac{n_a}{n_1} + o(1) \implies \delta_{\gamma n_a}(L^{\text{in}}) \leq \gamma \frac{n_a}{n_1} + o(1). \quad \blacksquare$$

644 *Remark A.3.* We note that the RIP is only meaningful for  $\delta_{\gamma n_a} < 1$ . Hence the above  
 645 theorem is only meaningful for  $\gamma < \frac{n_1}{n_a} - o(1)$ . To avoid this complicating technicality, we  
 646 henceforth assume that  $a = 1$ , *i.e.* that the target cluster is  $C_1$ .

### 647 **A.2. Bounding the size of the Perturbation.**

648 *Theorem A.4.* Suppose that  $G \sim \mathcal{G}_n$  with  $\mathcal{G}_n$  satisfying (A3). If  $L$  denotes the Laplacian  
 649 of  $G$  and  $M := L - L^{\text{in}}$  then  $\|M\|_2 \leq o(1)$ .

650 *Proof.* Letting  $\delta_{ij}$  denote the Kronecker delta symbol, observe that

$$651 \quad L_{ij} := \delta_{ij} - \frac{1}{d_i} A_{ij} = \delta_{ij} - \frac{1}{d_i^{\text{in}} + d_i^{\text{out}}} (A_{ij}^{\text{in}} + A_{ij}^{\text{out}}).$$

652 Earlier we defined  $r_i = d_i^{\text{out}}/d_i^{\text{in}}$ . We now use the following easily verifiable identity:

$$653 \quad \frac{1}{d_i^{\text{in}} + d_i^{\text{out}}} = \frac{1}{d_i^{\text{in}}} - \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right).$$

654 Thus:

$$\begin{aligned} 655 \quad L_{ij} &= \delta_{ij} - \left( \frac{1}{d_i^{\text{in}}} - \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) \right) (A_{ij}^{\text{in}} + A_{ij}^{\text{out}}) \\ 656 \quad &= \left( \delta_{ij} - \frac{1}{d_i^{\text{in}}} A_{ij}^{\text{in}} \right) - \frac{1}{d_i^{\text{in}}} A_{ij}^{\text{out}} + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) (A_{ij}^{\text{in}} + A_{ij}^{\text{out}}) \\ 657 \quad &= L_{ij}^{\text{in}} - \frac{1}{d_i^{\text{in}}} \left( 1 - \frac{r_i}{r_i + 1} \right) A_{ij}^{\text{out}} + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) A_{ij}^{\text{in}} \\ 658 \quad &= L_{ij}^{\text{in}} - \frac{1}{d_i^{\text{in}}} \left( \frac{1}{r_i + 1} \right) A_{ij}^{\text{out}} + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) A_{ij}^{\text{in}}. \end{aligned}$$

660 That is,  $M_{ij} = -\frac{1}{d_i^{\text{in}}} \left( \frac{1}{r_i + 1} \right) A_{ij}^{\text{out}} + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) A_{ij}^{\text{in}}$ . To bound the spectral norm we use Gersh-  
 661 gorin's disks, noting that  $M_{ii} = 0$  for all  $i$ :

$$\begin{aligned} 662 \quad \|M\|_2 &= \max_i \{ |\mu_i| : \mu_i \text{ eigenvalue of } M \} \leq \max_i \sum_j |M_{ij}| \\ 663 \quad &= \max_i \frac{1}{d_i^{\text{in}}} \left( \frac{1}{r_i + 1} \right) \sum_j A_{ij}^{\text{out}} + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) \sum_j A_{ij}^{\text{in}} \\ 664 \quad &= \max_i \left\{ \frac{1}{d_i^{\text{in}}} \left( \frac{1}{r_i + 1} \right) (d_i^{\text{out}}) + \frac{1}{d_i^{\text{in}}} \left( \frac{r_i}{r_i + 1} \right) (d_i^{\text{in}}) \right\} \\ 665 \quad &= \max_i \left\{ \left( \frac{r_i}{r_i + 1} \right) + \left( \frac{r_i}{r_i + 1} \right) \right\} \leq 2 \max_i r_i \leq 2\epsilon_2 = o(1) \\ 666 \end{aligned}$$

667 by (A3). ■

668 *Theorem A.5.* Suppose that  $G \sim \mathcal{G}_n$  with  $\mathcal{G}_n$  satisfying (A1)–(A4). If  $L$  denotes the Lapla-  
 669 cian of  $G$  and  $|C_1 \triangle \Omega| = \epsilon n_1$  with  $\epsilon \leq 0.13$  then  $\epsilon_{\mathbf{y}} = o(1)$  and  $\epsilon_L^{\gamma n_1} = o(1)$  for any  $\gamma \in (0, 1)$ .

670

671 *Proof.* Recall that  $\epsilon_{\mathbf{y}} = \frac{\|\mathbf{e}\|_2}{\|\mathbf{y}^{\text{in}}\|_2}$  and  $\epsilon_L^{\gamma_{n_1}} = \frac{\|M\|_2^{(\gamma_{n_1})}}{\|L^{\text{in}}\|_2^{(\gamma_{n_1})}}$ . Using the bound on the restricted  
672 isometry constant of  $L^{\text{in}}$  from Theorem A.2 we have:

$$\begin{aligned} 673 \quad \|\mathbf{y}^{\text{in}}\|_2^2 &= \|L^{\text{in}}(\mathbf{1}_W - \mathbf{1}_U)\|_2^2 \geq (1 - \delta_{\epsilon_{n_1}}(L^{\text{in}})) \|\mathbf{1}_W - \mathbf{1}_U\|_2^2 \\ 674 \quad &\geq (\epsilon - o(1)) |C_1 \triangle \Omega| = (\epsilon^2 - o(1)) n_1 \quad \blacksquare \end{aligned}$$

Thus  $\|\mathbf{y}^{\text{in}} + \mathbf{z}^{\text{in}}\|_2 \geq \sqrt{\epsilon^2 - o(1)} \sqrt{n_1}$ . On the other hand:

$$\|\mathbf{e}\|_2 = \|\mathbf{y} - \mathbf{y}^{\text{in}}\|_2 = \|L\mathbf{1}_\Omega - L^{\text{in}}\mathbf{1}_\Omega\|_2 = \|M\mathbf{1}_\Omega\|_2 \leq \|M\|_2 \|\mathbf{1}_\Omega\|_2 \leq o(1) \sqrt{(1 + \epsilon)n_1}$$

Thus:

$$\epsilon_{\mathbf{y}} = \frac{\|\mathbf{e}\|_2}{\|\mathbf{y}^{\text{in}}\|_2} \leq \frac{o(1) \sqrt{(1 + \epsilon)} \sqrt{n_1}}{\sqrt{(\epsilon^2 - o(1))} \sqrt{n_1}} = o(1)$$

as  $\epsilon$  is a constant, *i.e.* independent of  $n_1$ . The bound on  $\epsilon_L^{\gamma_{n_1}}$  is easier. By Lemma 2.6 and Property 3:

$$\|L^{\text{in}}\|_2^{(\gamma_{n_1})} \geq \sigma_{\gamma_{n_1-1}}(L^{\text{in}}) = \lambda_{\gamma_{n_1-1}}(L^{\text{in}}) \geq \lambda_{k+1}(L^{\text{in}})$$

as long as  $\gamma_{n_1} \geq k + 3$ , which is certainly the case for large enough  $n_1$ . Because  $\lambda_1(L_{G_{C_1}}) = \dots = \lambda_1(L_{G_{C_k}}) = 0$  and the spectrum of  $L^{\text{in}}$  is the union of the spectra of the  $L_{G_{C_a}}$ , it follows that:

$$\lambda_{k+1}(L^{\text{in}}) = \min_{a=1}^k \lambda_2(L_{G_{C_a}}) \geq 1 - \epsilon_1 = 1 - o(1)$$

by (A1). By Theorem A.4 and Lemma 2.6  $\|M\|_2^{(\gamma_{n_1})} \leq \|M\|_2 = o(1)$ . It follows that:

$$\epsilon_L^{\gamma_{n_1}} = \frac{\|M\|_2^{(\gamma_{n_1})}}{\|L^{\text{in}}\|_2^{(\gamma_{n_1})}} = \frac{o(1)}{1 - o(1)} = o(1).$$

676 **A.3. Restricted Isometry Property for  $L$ .** Finally, we extend from  $\delta_s(L^{\text{in}})$  to  $\delta_s(L)$  using  
677 the following result of Herman and Strohmer (cf. [29]):

678 **Theorem A.6.** *Suppose that  $\Phi = \hat{\Phi} + M$ . Let  $\hat{\delta}_s$  and  $\delta_s$  denote the  $s$  restricted isometry  
679 constants of  $\hat{\Phi}$  and  $\Phi$  respectively. Then:*

$$680 \quad \delta_s \leq (1 + \hat{\delta}_s)(1 + \epsilon_\Phi^s)^2 - 1.$$

681 **Corollary A.7.** *Let  $L$  denote the Laplacian of  $G \sim \mathcal{G}_n$  satisfying (A1)–(A4). Then we have  
682  $\delta_{\gamma_{n_1}}(L) \leq \gamma + o(1)$  for any  $\gamma \in (0, 1)$ .*

*Proof.* By Theorem A.6 we have that:

$$\delta_{\gamma_{n_1}}(L) \leq (1 + \delta_{\gamma_{n_1}}(L^{\text{in}}))(1 + \epsilon_L^{\gamma_{n_1}})^2 - 1.$$

683 Substituting the values of  $\delta_{\gamma_{n_1}}(L^{\text{in}})$  and  $\epsilon_L^{\gamma_{n_1}}$  from Theorems A.2 and A.5 yields the claim.  $\blacksquare$



684 **B. Proof of Theorem 4.1.** Before proving this theorem we prove the a series of lemmas.  
 685 We first note that Assumptions (A3) and (A4) easily allow us to bound  $\text{vol}(S)$ , which will be  
 686 required in the proof of Theorem 4.1:

687 **Lemma B.1.** *Suppose that  $\mathcal{G}_n$  satisfies (A3) and (A4). For any  $S \subset V$  define  $\text{vol}^{\text{in}}(S) =$   
 688  $\sum_i d_i^{\text{in}}$ . Then for any  $G \in \mathcal{G}_n$  we have that:*

689 (1)  $(1 - \epsilon_3)|S|d_{\text{av}}^{\text{in}} \leq \text{vol}^{\text{in}}(S) \leq (1 + \epsilon_3)|S|d_{\text{av}}^{\text{in}}$ ; and (2)  $\text{vol}^{\text{in}}(S) \leq \text{vol}(S) \leq (1 + \epsilon_2)\text{vol}^{\text{in}}(S)$ .

690 *Proof.* For part (1), observe that:

691 
$$\text{vol}^{\text{in}}(S) = \sum_{i \in S} d_i^{\text{in}} \geq |S|d_{\text{min}}^{\text{in}} \geq |S|(1 - \epsilon_3)d_{\text{av}}^{\text{in}},$$

where the final inequality is from (A4). The bound  $\text{vol}^{\text{in}}(S) \leq (1 + \epsilon_3)|S|d_{\text{av}}^{\text{in}}$  follows similarly.  
 For part (2) we note that by assumption (A3)  $d_i = d_i^{\text{in}} + d_i^{\text{out}} \leq d_i^{\text{in}} + \epsilon_2 d_i^{\text{in}} = (1 + \epsilon_2)d_i^{\text{in}}$ .  
 Hence:

$$\text{vol}(S) = \sum_{i \in S} d_i \leq \sum_{i \in S} (1 + \epsilon_2)d_i^{\text{in}} = (1 + \epsilon_2)\text{vol}^{\text{in}}(S)$$

692 while the lower bound follows simply from the fact that  $d_i \geq d_i^{\text{in}}$ . ■

**Lemma B.2.** *Let  $G \in \mathcal{G}_n$  satisfies Assumptions (A1)–(A4). If  $N_{G_{C_1}} := D_{G_{C_1}}^{-1/2} A_{G_{C_1}} D_{G_{C_1}}^{-1/2}$   
 and  $U, \Gamma \subset C_1$  then:*

$$\left| \langle D_{G_{C_1}}^{1/2} \mathbf{1}_U, N_{G_{C_1}}^t D_{G_{C_1}}^{1/2} \mathbf{1}_\Gamma \rangle - \frac{\text{vol}^{\text{in}}(U)\text{vol}^{\text{in}}(\Gamma)}{\text{vol}^{\text{in}}(G_{C_1})} \right| \leq \epsilon_1^t \sqrt{\text{vol}^{\text{in}}(U)\text{vol}^{\text{in}}(\Gamma)}$$

*Proof.* From the proof of Lemma 2 in [14] (note that they use  $M_{G_{C_1}}$  instead of  $N_{G_{C_1}}$ ) we  
 get that:

$$\left| \langle D_{G_{C_1}}^{1/2} \mathbf{1}_U, N_{G_{C_1}}^t D_{G_{C_1}}^{1/2} \mathbf{1}_\Gamma \rangle - \frac{\text{vol}^{\text{in}}(U)\text{vol}^{\text{in}}(\Gamma)}{\text{vol}^{\text{in}}(G_{C_1})} \right| \leq \lambda_{n_1-1}(N_{G_{C_1}})^t \sqrt{\text{vol}^{\text{in}}(U)\text{vol}^{\text{in}}(\Gamma)}$$

693 By Lemma 2.3 and (A2) we get that  $\lambda_{n_1-1}(N_{G_{C_1}}) = 1 - \lambda_2(L_{G_{C_1}}) \leq \epsilon_1$ . ■

694 *Proof of Theorem 4.1.* As in §3, let  $U = C_1 \setminus \Omega$  and  $W = \Omega \setminus C_1$ . Let  $|U| = un_1$ , in which  
 695 case  $|W| = (\epsilon + u)n_1$ . We shall prove that  $u = o(1)$ . By definition,  $\Omega$  is the set of the  $(1 + \epsilon)n_1$   
 696 largest entries in  $\mathbf{v}^{(t)} := P^t D \mathbf{1}_\Gamma$ . Because  $U$  is not in  $\Omega$ , but  $W$  is, we must have  $v_i^{(t)} \leq v_j^{(t)}$   
 697 for every  $i \in U$  and  $j \in W$ . We sum first over  $j \in W$  and then sum over  $i \in U$  to obtain:

698 
$$v_i^{(t)} \leq v_j^{(t)} \implies (\epsilon + u)n_1 v_i^{(t)} \leq \sum_{j \in W} v_j^{(t)} \implies (\epsilon + u)n_1 \sum_{i \in U} v_i^{(t)} \leq un_1 \sum_{j \in W} v_j^{(t)}.$$

699 It follows that:

700 (B.1) 
$$\sum_{i \in U} v_i^{(t)} \leq \frac{u}{\epsilon + u} \sum_{j \in W} v_j^{(t)} \leq \sum_{j \in W} v_j^{(t)}.$$

701 Looking ahead, we shall show that if inequality (B.1) holds then  $u = o(1)$ .

702

We first show that the term on the left-hand side of inequality B.1, *i.e.* the sum over the vertices in  $C_1$  that were missed by  $\Omega$ , is necessarily quite large. We do this by relating  $P$  to  $P^{\text{in}}$ , the random walk transition matrix for the graph  $G^{\text{in}}$ . Note that  $G^{\text{in}}$  is a disjoint union of the graphs  $G_{C_a}$ . For every  $i \in [n]$ , define  $q_i := d_i^{\text{in}}/d_i$ . Observe that  $1/d_i = q_i/d_i^{\text{in}}$  and thus  $D^{-1} = D_{\text{in}}^{-1}Q$  where  $Q$  is the diagonal matrix with  $(i, i)$ -th entry  $q_i$ . Now:

$$P = AD^{-1} = (A^{\text{in}} + A^{\text{out}})D^{-1} = A^{\text{in}}(D_{\text{in}}^{-1}Q) + A^{\text{out}}D^{-1} = P^{\text{in}}Q + A^{\text{out}}D^{-1}.$$

Observe that  $P$ ,  $P^{\text{in}}Q$  and  $A^{\text{out}}D^{-1}$  all have non-negative entries. It follows that for any non-negative vector  $\mathbf{x}$ :  $P\mathbf{x}$  and  $P^{\text{in}}Q\mathbf{x}$  are also non-negative and  $P\mathbf{x} \geq P^{\text{in}}Q\mathbf{x}$ , where the inequality should be interpreted componentwise. One can extend the inequality by iterated multiplication:

$$P^t\mathbf{x} \geq (P^{\text{in}}Q)^t\mathbf{x} \geq q_{\min}^t (P^{\text{in}})^t\mathbf{x}$$

703 and again the inequality should be interpreted componentwise. Now:

$$\begin{aligned} 704 \quad \sum_{i \in U} v_i^{(t)} &= \langle \mathbf{1}_U, \mathbf{v}^{(t)} \rangle = \langle \mathbf{1}_U, P^t D \mathbf{1}_\Gamma \rangle \geq \langle \mathbf{1}_U, q_{\min}^t (P^{\text{in}})^t D \mathbf{1}_\Gamma \rangle \\ 705 \quad &= q_{\min}^t \langle \mathbf{1}_U, (P^{\text{in}})^t D_{\text{in}} \mathbf{1}_\Gamma \rangle = q_{\min}^t \langle \mathbf{1}_U, (P_{G_{C_1}})^t D_{G_{C_1}} \mathbf{1}_\Gamma \rangle, \\ 706 \end{aligned}$$

707 where the final line follows as  $U, \Gamma \subset C_1$ .

708 Our goal now is to bound the quantity  $\langle \mathbf{1}_U, (P_{G_{C_1}})^t D_{G_{C_1}} \mathbf{1}_\Gamma \rangle$ . One can rearrange the  
709 iterated matrix product slightly:

$$\begin{aligned} 710 \quad (P_{G_{C_1}})^t &= (A_{G_{C_1}} D_{G_{C_1}}^{-1})^t = A_{G_{C_1}} D_{G_{C_1}}^{-1} A_{G_{C_1}} D_{G_{C_1}}^{-1} \dots A_{G_{C_1}} D_{G_{C_1}}^{-1} \\ 711 \quad &= D_{G_{C_1}}^{1/2} (D_{G_{C_1}}^{-1/2} A_{G_{C_1}} D_{G_{C_1}}^{-1/2}) (D_{G_{C_1}}^{-1/2} A_{G_{C_1}} D_{G_{C_1}}^{-1/2}) \dots (D_{G_{C_1}}^{-1/2} A_{G_{C_1}} D_{G_{C_1}}^{-1/2}) D_{G_{C_1}}^{-1/2} \\ 712 \quad &= D_{G_{C_1}}^{1/2} N_{G_{C_1}}^t D_{G_{C_1}}^{-1/2}, \\ 713 \end{aligned}$$

714 Hence, we have

$$\begin{aligned} 715 \quad \langle \mathbf{1}_U, (P_{G_{C_1}})^t D_{G_{C_1}} \mathbf{1}_\Gamma \rangle &= \langle \mathbf{1}_U, (D_{G_{C_1}}^{1/2} N_{G_{C_1}}^t D_{G_{C_1}}^{-1/2}) D_{G_{C_1}} \mathbf{1}_\Gamma \rangle \\ 716 \quad &= \langle D_{G_{C_1}}^{1/2} \mathbf{1}_U, N_{G_{C_1}}^t D_{G_{C_1}}^{1/2} \mathbf{1}_\Gamma \rangle \geq \frac{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)}{\text{vol}^{\text{in}}(G_{C_1})} - \epsilon_1^t \sqrt{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)}, \\ 717 \end{aligned}$$

718 where the final inequality follows from Lemma B.2. Returning to (20):

$$719 \quad (\text{B.2}) \quad \sum_{i \in U} v_i^{(t)} \geq q_{\min}^t \left( \frac{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)}{\text{vol}^{\text{in}}(G_{C_1})} - \epsilon_1^t \sqrt{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)} \right).$$

We now consider the right hand side of (B.1), *i.e.* the sum over  $W$ . Because  $W \subset V \setminus C_1$  we have that:

$$\sum_{j \in W} v_j^{(t)} \leq \sum_{j \in V \setminus C_1} |v_j^{(t)}| = \|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1$$

Thus it remains to bound  $\|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1$ . Observe that:

$$\mathbf{v}_{V \setminus C_1}^{(t)} = A^{\text{in}} D^{-1} \mathbf{v}_{V \setminus C_1}^{(t-1)} + \left( A^{\text{out}} D^{-1} \mathbf{v}^{(t-1)} \right)_{V \setminus C_1}.$$

Clearly

$$\left\| \left( A^{\text{out}} D^{-1} \mathbf{v}^{(t-1)} \right)_{V \setminus C_1} \right\|_1 \leq \left\| A^{\text{out}} D^{-1} \mathbf{v}^{(t-1)} \right\|_1$$

and so  $\|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1 \leq$

$$\|A^{\text{in}} D^{-1} \mathbf{v}_{V \setminus C_1}^{(t-1)}\|_1 + \|A^{\text{out}} D^{-1} \mathbf{v}^{(t-1)}\|_1 \leq \|A^{\text{in}} D^{-1}\|_1 \|\mathbf{v}_{V \setminus C_1}^{(t-1)}\|_1 + \|A^{\text{out}} D^{-1}\|_1 \|\mathbf{v}^{(t-1)}\|_1$$

Moreover:  $\|A^{\text{in}} D^{-1}\|_1 = \max_j \sum_i \frac{A_{ij}^{\text{in}}}{d_j} = \max_j \frac{d_j^{\text{in}}}{d_j} \leq 1$  and similarly  $\|A^{\text{out}} D^{-1}\|_1 = \max_j \frac{d_j^{\text{out}}}{d_j} \leq \max_j r_j \leq \epsilon_2$  by assumption (A2). Thus  $\|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1 \leq 1 \|\mathbf{v}_{V \setminus C_1}^{(t-1)}\|_1 + \epsilon_2 \|\mathbf{v}^{(t-1)}\|_1$ . Solving this recursion relation we obtain:

$$\|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1 \leq \epsilon_2 \sum_{s=0}^{t-1} \|\mathbf{v}^{(s)}\|_1 + \|\mathbf{v}_{V \setminus C_1}^{(0)}\|_1$$

720 Because  $\mathbf{v}^{(0)} = D \mathbf{1}_\Gamma$  and  $\Gamma \subset C_1$ , it follows that  $\|\mathbf{v}_{V \setminus C_1}^{(0)}\|_1 = 0$  and  $\|\mathbf{v}^{(0)}\|_1 = \text{vol}(\Gamma)$ . Because  
721  $\|P\|_1 = 1$  it follows that  $\|\mathbf{v}^{(s)}\|_1 = \|\mathbf{v}^{(0)}\|_1 = \text{vol}(\Gamma)$  for all  $s$ . Thus:

$$722 \quad (\text{B.3}) \quad \sum_{j \in W} v_j^{(t)} \leq \|\mathbf{v}_{V \setminus C_1}^{(t)}\|_1 \leq t \epsilon_2 \text{vol}(\Gamma) \leq t \epsilon_2 (1 + \epsilon_2) \text{vol}^{\text{in}}(\Gamma),$$

723 where the final inequality follows from Lemma B.1. Now let us put this all together. Returning  
724 to (B.1) with (B.2) and (B.3) in hand:

$$725 \quad (\text{B.4}) \quad q_{\min}^t \left( \frac{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)}{\text{vol}^{\text{in}}(G_{C_1})} - \epsilon_1^t \sqrt{\text{vol}^{\text{in}}(U) \text{vol}^{\text{in}}(\Gamma)} \right) \leq t \epsilon_2 (1 + \epsilon_2) \text{vol}^{\text{in}}(\Gamma)$$

$$726 \quad (\text{B.5}) \quad \implies q_{\min}^t \left( \frac{\text{vol}^{\text{in}}(U)}{\text{vol}^{\text{in}}(G_{C_1})} - \epsilon_1^t \sqrt{\frac{\text{vol}^{\text{in}}(U)}{\text{vol}^{\text{in}}(\Gamma)}} \right) \leq t \epsilon_2 (1 + \epsilon_2).$$

727

728 From Lemma B.1 and the assumptions on  $|U|$  and  $|\Gamma|$ :

$$729 \quad \frac{\text{vol}^{\text{in}}(U)}{\text{vol}^{\text{in}}(G_{C_1})} \geq \frac{(1 - \epsilon_3) d_{\text{av}}^{\text{in}} |U|}{(1 + \epsilon_3) d_{\text{av}}^{\text{in}} |C_1|} = \frac{(1 - \epsilon_3) u n_1}{(1 + \epsilon_3) n_1} = \frac{1 - \epsilon_3}{1 + \epsilon_3} u$$

$$730 \quad \frac{\text{vol}^{\text{in}}(U)}{\text{vol}^{\text{in}}(\Gamma)} \leq \frac{(1 + \epsilon_3) d_{\text{av}}^{\text{in}} |U|}{(1 - \epsilon_3) d_{\text{av}}^{\text{in}} |\Gamma|} \leq \frac{(1 + \epsilon_3)}{(1 - \epsilon_3)} \frac{u}{g \epsilon_1^{2t-1}}$$

731

Finally because  $q_i = 1 - r_i$  it follows that  $q_{\min} \geq 1 - \epsilon_2$ . Putting this all into equation (B.5):

$$(1 - \epsilon_2)^t \left( \frac{1 - \epsilon_3}{1 + \epsilon_3} u - \epsilon_1^{1/2} \sqrt{\frac{(1 + \epsilon_3) u}{(1 - \epsilon_3) g}} \right) \leq t \epsilon_2 (1 + \epsilon_2)$$

At this stage it is illuminating to use the assumption that  $\epsilon_1, \epsilon_2, \epsilon_3 = o(1)$ . Observe that:

$$\frac{1 - o(1)}{1 + o(1)} = 1 - o(1), \quad \frac{1 + o(1)}{1 - o(1)} = 1 + o(1), \quad \text{and } (1 - o(1))^t = 1 - o(1)$$

where the final equality follows as  $t$  is constant with respect to  $n$ . Hence:

$$(1 - o(1))u - o(\sqrt{u}) \leq o(1) \implies u \leq o(1) + o(u).$$

732 This is only possible if  $u = o(1)$ . It follows that  $|C_1 \triangle \Omega| = |U| + |W| = (\epsilon + 2u)n_1 = (\epsilon + o(1))n_1$   
733 as stated. ■

734 **C. Showing the SBM satisfies our assumptions.** Let us verify that  $\text{SBM}(\mathbf{n}, P)$  satisfies  
735 the assumptions (A1)–(A4), under the hypotheses of Theorem 9.3. Recall that our assumption  
736 is  $P_{ab} = (\beta + o(1)) \log(n)/n$  for  $a \neq b$ , and that  $P_{aa} = \omega \log(n)/n_a$  for  $a = 1, \dots, k$ . As we also  
737 assume that  $n_1 = O(n) \rightarrow \infty$ , and  $n_1$  is the size of the smallest cluster, we get that  $k = O(1)$ ,  
738 *i.e.* (A1) holds.

739 **Theorem C.1** (see [7, 8]). *Let  $G \sim \text{ER}(n, q)$  with  $q = (\beta + o(1)) \log(n)/n$ . There exist a*  
740 *function  $\eta(\beta)$  satisfying  $0 < \eta(\beta) < 1$  and  $\lim_{\beta \rightarrow \infty} \eta(\beta) = 0$  such that*

$$741 \quad d_{\max}(G) = (1 + \eta(\beta))\beta \log n + o(1) \leq 2\beta \log(n) + o(1) \text{ a.s.}$$

742 **Theorem C.2** (see [23], Theorem 3.4 (ii)). *If  $G \sim \text{ER}(n_a, p)$  with  $p_a = \omega \log(n)/n_a$  where*  
743  *$\omega \rightarrow \infty$ , then  $d_{\min}(G) = (1 - o(1))\omega \log(n)$  and  $d_{\max}(G) = (1 + o(1))\omega \log(n)$  a.s.*

744 **Theorem C.3.** *Suppose that  $G \sim \text{ER}(n_a, p)$  with  $p = \omega \log(n)/n_a$  where  $\omega \rightarrow \infty$ . Then we*  
745 *have almost surely  $|\lambda_i(L) - 1| = O(\omega^{-1/2}) = o(1)$  for all  $i > 1$ .*

*Proof.* Theorem 4 in [15] shows that

$$|\lambda_i(L^{\text{sym}}) - 1| \leq \sqrt{\frac{6 \log(2n_a)}{\omega \log(n)}}.$$

746 By Lemma 2.3  $L^{\text{sym}}$  and  $L$  have the same spectrum. The result follows as  $\log(n) \geq \log(n_a)$  ■

747 As each  $G_{C_a} \sim \text{ER}(n_a, p)$ , it follows from Theorem C.3 that:

748 **Corollary C.4.** *SBM( $\mathbf{n}, P$ ) with parameters as in Theorem 9.3 satisfies assumption (A2)*  
749 *with  $\epsilon_1 = O(\omega^{-1/2})$ .*

750 We now discuss the remaining two assumptions. Let  $G^{\text{in}}$  and  $G^{\text{out}}$  be as in §2. If  $G \sim$   
751  $\text{SBM}(\mathbf{n}, P)$  then  $G^{\text{in}}$  consists of  $k$  disjoint Erdős - R enyi graphs,  $G_{C_a} \sim \text{ER}(n_a, p)$ . The graph  
752  $G^{\text{out}}$  is not an Erdős - R enyi graph, as there is zero probability of it containing an edge between  
753 two vertices in the same cluster (because we have removed them). However, we can profitably  
754 think of  $G^{\text{out}}$  as a subgraph of some  $\widetilde{G}^{\text{out}} \sim \text{ER}(n, q)$ . In particular, any upper bounds on the  
755 degrees of vertices in  $\widetilde{G}^{\text{out}}$  are automatically bounds on the degrees in  $G^{\text{out}}$ . Thus, we have  
756 the following corollaries of Theorems C.2 and C.1:

757 **Corollary C.5.** *If  $G \sim \text{SBM}(\mathbf{n}, P)$  with parameters as in Theorem 9.3 then  $d_{\max}^{\text{out}}(G) \leq$*   
758  *$2\beta \log n + o(1)$  a.s.*

759 *Proof.* Consider  $G^{\text{out}}$  as a subgraph of  $\widetilde{G}^{\text{out}} \sim \text{ER}(n, q)$  and apply Theorem C.1 ■

760 **Corollary C.6.** *If  $G \sim \text{SBM}(\mathbf{n}, P)$  with parameters as in Theorem 9.3, then  $d_{\min}^{\text{in}}(G) \geq$*   
 761  *$(1 - o(1))\omega \log(n)$  and  $d_{\max}^{\text{in}}(G) \leq (1 + o(1))\omega \log(n)$  a.s.*

762 *Proof.* If  $i \in C_a$  then  $d_i^{\text{in}} = d_i(G_{C_a})$ , where  $G_{C_a} \sim \text{ER}(n_a, p)$ . Clearly:

$$763 \quad d_{\max}^{\text{in}}(G) = \max_i d_i^{\text{in}} = \max_a d_{\max}(G_{C_a})$$

764 By Theorem C.2,  $d_{\max}(G_a) = (1 + o(1))\omega \log(n)$  a.s. Note that the  $d_{\max}(G_{C_a})$  are independent  
 765 random variables, and since we are taking a maximum over  $k = \mathcal{O}(1)$  of them, it follows that  
 766  $\max_a d_{\max}(G_{C_a}) \leq (1 + o(1))\omega \log(n)$  a.s. too. The proof for  $d_{\min}^{\text{in}}(G)$  is similar. ■

767 **Corollary C.7.**  *$\text{SBM}(\mathbf{n}, P)$  with parameters as in Theorem 9.3 satisfies assumption (A3)*  
 768 *with  $\epsilon_2 = O(\omega^{-1})$ .*

769 *Proof.* First of all, it is clear that for any  $i$ ,  $d_i^{\text{out}}/d_i^{\text{in}} \leq d_{\max}^{\text{out}}/d_{\min}^{\text{in}}$ . From Corollaries C.5  
 770 and C.6 we have:

$$771 \quad \frac{d_{\max}^{\text{out}}}{d_{\min}^{\text{in}}} \leq \frac{2\beta \log n + o(1)}{(1 - o(1))\omega \log(n)} = \frac{2\beta + o(1)}{(1 - o(1))\omega} = O(\omega^{-1}). \quad \blacksquare$$

773 **Corollary C.8.**  *$\text{SBM}(\mathbf{n}, P)$  with parameters as in Theorem 9.3 satisfies assumption (A4).*

774 *Proof.* Observe that  $d_{\text{av}}^{\text{in}} = \omega \log(n)$ . The result then follows from Corollary C.6. ■