

# A sharp PageRank algorithm with applications to edge ranking and graph sparsification

Fan Chung\* and Wenbo Zhao

University of California, San Diego  
La Jolla, CA 92093  
{fan,w3zhao}@ucsd.edu

**Abstract.** We give an improved algorithm for computing personalized PageRank vectors with tight error bounds which can be as small as  $O(n^{-k})$  for any fixed positive integer  $k$ . The improved PageRank algorithm is crucial for computing a quantitative ranking for edges in a given graph. We will use the edge ranking to examine two interrelated problems — graph sparsification and graph partitioning. We can combine the graph sparsification and the partitioning algorithms using PageRank vectors to derive an improved partitioning algorithm.

## 1 Introduction

PageRank, which was first introduced by Brin and Page [11], is at the heart of Google’s web searching algorithms. Originally, PageRank was defined for the Web graph (which has all webpages as vertices and hyperlinks as edges). For any given graph, PageRank is well-defined and can be used for capturing quantitative correlations between pairs of vertices as well as pairs of subsets of vertices. In addition, PageRank vectors can be efficiently computed and approximated (see [3, 4, 10, 22, 26]). The running time of the approximation algorithm in [3] for computing a PageRank vector within an error bound of  $\epsilon$  is basically  $O(1/\epsilon)$ . For the problems that we will examine in this paper, it is quite crucial to have a sharper error bound for PageRank. In Section 2, we will give an improved approximation algorithm with running time  $O(m \log(1/\epsilon))$  to compute PageRank vectors within an error bound of  $\epsilon$ .

The PageRank algorithm is originally meant for determining the “importance” of vertices in the Web graph. It is also essential to identify the “importance” of edges in dealing with various problems. We will use PageRank vectors to define a qualitative ranking of edges that we call *Green* values for edges because of the connection with discrete Green functions. The Green values for edges can also be viewed as a generalized version of effective resistances in electric network theory. The detailed definition for Green values of edges will be given in Section 3. We then use the sharp approximate PageRank algorithm to compute Green values within sharp error bounds.

To illustrate the usage of Green values, we examine a basic problem on sparsifying graphs. Graph sparsification was first introduced by Benczúr and Karger [9, 30–32] for approximately solving various network design problems. The heart of the graph sparsification algorithms are the sampling techniques for randomly selecting edges. The goal is to approximate a given graph  $G$  with  $m$  edges on  $n$  vertices by a sparse graph  $\tilde{G}$ , called *sparsifier*, with  $O(n \log n)$  edges (or fewer) on the same set of vertices in such a way that every cut in sparsifier  $\tilde{G}$  is within a factor of  $(1 \pm \epsilon)$  of the corresponding cut in  $G$  for some constant  $\epsilon > 0$ . It was shown that, for all  $x \in \{0, 1\}^n$ ,

$$|x^T \tilde{L}x - x^T Lx| \leq \epsilon x^T Lx, \tag{1}$$

where  $L$  and  $\tilde{L}$  are the Laplacian of the graph  $G$  and its sparsifier  $\tilde{G}$ , respectively.

Spielman and Teng [46, 47] devised a sampling scheme to construct a *spectral sparsifier* with  $O(n \log^c n)$  edges for some (large) constant  $c$  in  $O(m \text{ polylog}(n))$  time. A spectral sparsifier  $\tilde{G}$  for graph  $G$  is a sparsifier satisfying equation (1) for all  $x \in \mathbb{R}^n$ . In [48] Spielman and Srivastava gave a different sampling

\* Research supported in part by ONR MURI N000140810747, and AF/SUB 552082

scheme using the effective resistance of electrical networks to construct an improved spectral sparsifier with  $O(n \log n)$  edges. In the process for constructing this spectral sparsifier, they need to use the Spielman-Teng solver [47] as subroutines for solving  $O(\log n)$  linear systems. The running time of the sparsification algorithm is mainly dominated by the running time of Spielman-Teng solver which is  $O(m \log^c n)$  for a very large constant  $c$  [46, 47]. Recently, Batson, Spielman and Srivastava [8] gave an elegant construction for a sparsifier with a linear number of edges although the running time is  $O(n^3 m)$ .

Here, we use Green values to sample edges of  $G$  in order to form the sparsifier  $\tilde{G}$  with  $O(n \log n)$  edges. There are two advantages of sampling using PageRank and Green values. The running time for our sparsification algorithm is significantly faster and simpler than those in [8, 47] since we avoid using Spielman-Teng solver for solving linear system. In addition, the graph sparsification problem is closely related to graph partitioning algorithms.

For graph partitioning algorithms, previously widely used approach is the recursive spectral method which finds a balanced cut in a graph on  $n$  vertices with running time  $O(n^2 \lambda^{-1} \text{polylog}(n))$  (see [45]), together with an approximation guarantee within a quadratic root of the optimal conductance (where  $\lambda$  denotes the spectral gap of the normalized Laplacian). The running time can be further improved to  $O(n^2 \text{polylog}(n))$  by using Spielman-Teng solver for linear systems [47]. Another approach for the balanced cut problem is by using commodity flows [7, 41]. In [7] the approximation guarantee is within a factor of  $\log n$  of the optimal in [7]), which was further reduced to  $\log^{1/2} n$  in [6] but the running time is still  $O(n^2 \text{polylog}(n))$ . In another direction, Spielman and Teng [46, 48] introduced local graph partitioning which yields a cut near the specified seeds with running time only depending on the volume of the output. Their local partitioning algorithm has an approximation guarantee similar to the spectral method by using a mixing result on random walks [38, 39]. Andersen, Chung and Lang [3] used PageRank vectors to give a local partitioning algorithm with improved approximation guarantee and running time. Recently, Andersen and Peres use involving sets instead of PageRank to further improved the running time [4].

Our balanced-cut algorithm consisting of two parts. First we use PageRank vectors to sparsify the graph. Then we use the known PageRank partitioning algorithm to find a balanced cut. Both parts have the same complexity as computing the PageRank vectors. Consequently, the complexity for our PageRank balanced-cut algorithm is  $O(m \log^2 n / \phi + n \text{polylog}(n))$  for any input graph on  $n$  vertices and  $m$  edges. The balanced-cut algorithm here can be viewed as an application of graph sparsification.

## 2 A sharp PageRank approximation algorithm

We consider an undirected, weighted graph  $G = (V, E, w)$  with  $n$  vertices and  $m$  edges where the edge weights  $w(u, v) = w(v, u) \geq 0$  and the edge set  $E$  consists of all pairs  $(u, v)$  with  $w(u, v) > 0$ . The weighted degree  $d(u)$  of vertex  $u$  is the sum of  $w(u, v)$  over all  $v$ , i.e.,  $d(u) = \sum_v w(u, v)$ .

A typical random walk is defined by its transition probability matrix  $P$  satisfying  $P(u, v) = w(u, v)/d(u)$ . We may write  $P = D^{-1}A$ , where  $A$  is the weighted adjacency matrix satisfying  $A(u, v) = w(u, v)$  for all pairs of  $u, v \in V$  and  $D$  is the diagonal matrix of weighted degree. We here consider the *lazy walk*  $Z$  on  $G$ , defined by

$$Z = \frac{I + P}{2}.$$

In [3] PageRank vector  $p$  satisfies a recurrence relation involving a seed vector  $s$  (as a probability distribution) and a positive jumping constant  $\alpha < 1$  (or transportation constant):

$$p = \alpha s + (1 - \alpha)pZ$$

where  $p$  and  $s$  are taken to be row vectors. In this paper, we consider the PageRank  $p$  as a discrete Green's function satisfying

$$\begin{aligned} p &= \alpha s (I - (1 - \alpha)Z)^{-1} \\ &= \beta s (\beta I + \mathbf{L})^{-1} \end{aligned}$$

where  $\beta = 2\alpha/(1 - \alpha)$  and  $\mathbf{L} = I - P$ . Note that the usual Green's function is associated with the pseudo inverse of  $L$ . Another way to express the recurrence of PageRank in terms of  $\beta$  is the following: For a positive value  $\beta > 0$ , the (personalized) PageRank vector  $\text{pr}_{\beta,s}$  with a seed vector  $s$  is the unique solution of the following:

$$\text{pr}_{\beta,s} = \frac{\beta}{2 + \beta} s + \frac{2}{2 + \beta} \text{pr}_{\beta,s} Z.$$

If a seed vector is the characteristic function  $\chi_u$  of a single vertex  $u$ , we may write  $\text{pr}_{\beta,\chi_u} = \text{pr}_{\beta,u}$  if there is no confusion. It is easy to check that  $\sum_{v \in V} \text{pr}_{\beta,s}(v) = 1$  since  $\sum_{v \in V} s(v) = 1$ .

The PageRank approximation algorithms in [3] contains the following steps, called **push** and **ApproximatePR**, which serve as subroutines later in the sharp approximate PageRank algorithm. For a vertex  $u$ , a current approximate PageRank vector  $p$  and a residual vector  $r$ , the **push** step is as follows:

**push**( $u$ ):

Let  $p' = p$  and  $r' = r$ , except for these changes:

1.  $p'(u) = p(u) + \frac{\beta}{2+\beta} r(u)$
2.  $r'(u) = \frac{r(u)}{2+\beta}$
3. For each vertex  $v$  such that  $(u, v) \in E$ :  
 $r'(v) = r(v) + \frac{r(u)}{(2+\beta)d(u)}$ .

**Lemma 1 ([3]).** *Let  $p'$  and  $r'$  denote the resulting vectors after performing **push**( $u$ ) on  $p$  and  $r$ . Then,*

$$p = \text{pr}_{\beta,s-r} \implies p' = \text{pr}_{\beta,s-r'}.$$

**ApproximatePR**( $s, \beta, \epsilon$ ):

1. Let  $p = \mathbf{0}$  and  $r = s$ .
2. While  $r(u) \geq \epsilon d(u)$  for some vertex  $u$ :
  - (a) Pick any vertex  $u$  where  $r(u) \geq \epsilon d(u)$ .
  - (b) Apply **push**( $u$ ).
3. Return  $p$  and  $r$ .

**Theorem 1 ([3]).** *For any  $s$  with  $\|s\|_1 \leq 1$ , and any constant  $\epsilon \in (0, 1]$ , the algorithm **ApproximatePR**( $s, \beta, \epsilon$ ) computes an approximate PageRank vector  $p = \text{pr}_{\beta,s-r}$  such that the residual vector  $r$  satisfies  $|r(v)/d(v)| \leq \epsilon$  for all  $v \in V$ . The running time of the algorithm is  $O(\frac{2+\beta}{\epsilon\beta})$ .*

We will improve the estimate error bound for the above algorithm by the following iterated process:

**SharpApproximatePR**( $s, \beta, \epsilon$ ):

1. Let  $\epsilon' = 1$ ,  $r = s$  and  $p = \mathbf{0}$ .
2. While  $\epsilon' > \epsilon$  :
  - (a) Set  $\epsilon' = \epsilon'/2$ .
  - (b) Let  $p'$  and  $r'$  be the output of algorithm **ApproximatePR**( $r, \beta, \epsilon'$ ).
  - (c) Let  $p = p + p'$  and  $r = r'$ .
3. Return  $p$  and  $r$ .

**Theorem 2.** For any constant  $\epsilon \in (0, 1]$ ,  $\beta > 0$  and any seed vector  $s$ , to approximate  $\text{pr}_{\beta, s}$  the algorithm **SharpApproximatePR**( $s, \beta, \epsilon$ ) computes approximate PageRank vector  $p = \text{pr}_{\beta, s-r}$  such that the residual vector  $r$  satisfies  $|r(v)/d(v)| \leq \epsilon$  for all  $v \in V$ . The running time of the algorithm is  $O(\frac{(2+\beta)m \log(1/\epsilon)}{\beta})$ .

*Proof.* Clearly, the algorithm returns an approximate PageRank vector  $p = \text{pr}_{\beta, s-r}$  where the residual vector satisfies that  $r(v) \leq \epsilon d(v)$  at every vertex.

To bound the running time, we consider one fixed round of the second step. Let  $T$  denote the total number of push operations performed by **ApproximatePR**( $r, \alpha, \epsilon'$ ), and let  $d_i$  be the degree of vertex involved in the  $i$ th push operation. When the  $i$ th push operation was performed, the PageRank at this vertex was at least  $\epsilon' d_i$ , so  $\|r\|_1$  decreased by at least

$$\frac{\beta \epsilon' d_i}{2 + \beta}.$$

Since  $\|r\|_1$  is at most  $2\epsilon' \sum_{v \in V} d(v) = 2m\epsilon'$ , we have

$$\frac{\beta \epsilon'}{2 + \beta} \sum_{i=1}^T d_i \leq 2m\epsilon',$$

so

$$\sum_{i=1}^T d_i \leq \frac{2(2 + \beta)m}{\beta}.$$

Note that there are at most  $\log(1/\epsilon)$  rounds in the second step. Thus, The total running time is bounded by  $O(\frac{(2+\beta)m \log(1/\epsilon)}{\beta})$ .  $\square$

As an immediate consequence of Theorem 2, we have the following by taking  $\epsilon$  to be the inverse of a power of  $n$ .

**Corollary 1.** In a graph on  $n$  vertices, for any fixed integer  $k$ ,  $\beta > 0$  and any seed vector  $s$ , to approximate  $\text{pr}_{\beta, s}$  the algorithm **SharpApproximatePR**( $s, \beta, \epsilon$ ) computes approximate PageRank vector  $p = \text{pr}_{\beta, s-r}$  such that the residual vector  $r$  satisfies  $|r(v)/d(v)| \leq n^{-k}$  for all  $v \in V$ . The running time of the algorithm is  $O(\frac{(2+\beta)m \log n}{\beta})$ .

### 3 The Green values for edges in a graph

Recall that graph  $G = (V, E, w)$  is a connected weighted undirected graph with  $n$  vertices,  $m$  edges. The combinatorial Laplacian of  $G$  is defined by  $L = D - A$  where  $A$  is the weighted adjacency matrix and  $D$  is the diagonal matrix of weighted degrees. If we orient the edges of  $G$  in an arbitrary but fixed way, we can write its Laplacian as  $L = B^T W B$ , where  $B_{m \times n}$  is the signed edge-vertex incidence matrix, given by

$$B(e, v) = \begin{cases} 1 & \text{if } v \text{ is } e\text{'s head,} \\ -1 & \text{if } v \text{ is } e\text{'s tail,} \\ 0 & \text{otherwise} \end{cases}$$

and  $W_{m \times m}$  is the diagonal matrix with  $W(e, e) = w(e)$ . The normalized Laplacian of  $G$  is defined to be  $\mathcal{L} = D^{-1/2} L D^{-1/2}$  and we write  $\mathcal{L} = S^T W S$  where  $S_{m \times n} = B D^{-1/2}$ . Since  $\mathcal{L}$  is symmetric and we have

$$\mathcal{L} = \sum_{i=0}^{n-1} \lambda_i \phi_i^T \phi_i,$$

where  $\lambda_0 = 0$  and  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{n-1} \leq 2$  are the nonzero eigenvalues of  $\mathcal{L}$  and  $\phi_0, \dots, \phi_{n-1}$  are a corresponding orthonormal basis of eigenvectors. Various properties concerning eigenvalues of the normalized Laplacian can be found in [13].

Denote the  $\beta$ -normalized Laplacian  $\mathcal{L}_\beta$  by  $\beta I + \mathcal{L}$ , and we may write  $\mathcal{L}_\beta = S'^T W_\beta S'$  where we define  $S'$  and  $W_\beta$  as follows:

$$S' = \begin{bmatrix} I \\ S \end{bmatrix}_{(n+m) \times n} \quad \text{and} \quad W_\beta = \begin{bmatrix} \beta I & 0 \\ 0 & W \end{bmatrix}_{(n+m) \times (n+m)}.$$

Here the index set for the columns of  $S'$  and columns (rows) of  $W_\beta$  is  $V \cup E$  where the first  $n$  columns (rows) are indexed by  $V$  and the last  $m$  columns (rows) are indexed by  $E$ .

Green's functions were first introduced in a celebrated essay by George Green [20] in 1828. Since then, the concept of Green's functions has been used in numerous areas, especially in the study of partial differential equations and quantum field theory. The discrete analog of Green's function which are associated with the normalized Laplacian of graphs were considered in [14] in connection with the study of Dirichlet eigenvalues with boundary conditions. In particular, the following modified Green's function  $\mathcal{G}_\beta$  was used in [14]. For  $\beta \in \mathbb{R}^+$ , let Green's function  $\mathcal{G}_\beta$  denote the symmetric matrix satisfying  $\mathcal{L}_\beta \mathcal{G}_\beta = I$ . Clearly, we have

$$\mathcal{G}_\beta = \sum_{i=0}^{n-1} \frac{1}{\lambda_i + \beta} \phi_i^T \phi_i. \quad (2)$$

We remark that the discrete Green's function is basically a symmetric form of the PageRank. Namely, it is straightforward to check that

$$\frac{\text{pr}_{\beta,s}}{\beta} = sD^{-1/2} \mathcal{G}_\beta D^{1/2}. \quad (3)$$

For each edge  $e = \{u, v\} \in E$ , we define the *Green value*  $g_\beta(u, v)$  of  $e$  to be a combination of four terms in PageRank vectors as follows:

$$\begin{aligned} g_\beta(u, v) &= \beta(\chi_u - \chi_v) D^{-1/2} \mathcal{G}_\beta D^{-1/2} (\chi_u - \chi_v)^T \\ &= \frac{\text{pr}_{\beta,u}(u)}{d(u)} - \frac{\text{pr}_{\beta,u}(v)}{d(v)} + \frac{\text{pr}_{\beta,v}(v)}{d(v)} - \frac{\text{pr}_{\beta,v}(u)}{d(u)}. \end{aligned} \quad (4)$$

From (2) and (4), we have the following fact which will be useful later.

**Lemma 2.** *The Green value  $g_\beta$  can be expressed by*

$$g_\beta(u, v) = \sum_{i=0}^{n-1} \frac{\beta}{\lambda_i + \beta} \left( \frac{\phi_i(u)}{\sqrt{d(u)}} - \frac{\phi_i(v)}{\sqrt{d(v)}} \right)^2.$$

**Lemma 3.** *For two distinct vertices  $u, v \in V$ , we have*

$$\frac{\beta}{2 + \beta} \left( \frac{1}{d(u)} + \frac{1}{d(v)} \right) \leq g_\beta(u, v) \leq \frac{1}{d(u)} + \frac{1}{d(v)}.$$

*Proof.* Since  $\lambda_i \leq 2$ , we have

$$\frac{\beta}{2 + \beta} \leq \frac{\beta}{\lambda_i + \beta} \leq 1.$$

From Lemma 2 we have

$$\frac{\beta}{2 + \beta} \left( \frac{\phi_i(u)}{\sqrt{d(u)}} - \frac{\phi_i(v)}{\sqrt{d(v)}} \right)^2 \leq g_\beta(u, v) \leq \sum_{i=0}^{n-1} \left( \frac{\phi_i(u)}{\sqrt{d(u)}} - \frac{\phi_i(v)}{\sqrt{d(v)}} \right)^2.$$

Note that for two fixed vertices  $u$  and  $v$ , the vector  $f_u$ , defined by  $f_u(i) = \phi_i(u)$ , for  $i = 0, 1, \dots, n-1$ , is orthogonal to  $f_v$ . This implies

$$\sum_{i=0}^{n-1} \frac{\phi_i(u)\phi_i(v)}{\sqrt{d(u)d(v)}} = 0$$

and

$$\sum_{i=0}^{n-1} \left( \frac{\phi_i(u)}{\sqrt{d(u)}} - \frac{\phi_i(v)}{\sqrt{d(v)}} \right)^2 = \sum_{i=0}^{n-1} \left( \frac{\phi_i^2(u)}{d(u)} + \frac{\phi_i^2(v)}{d(v)} \right) = \frac{1}{d(u)} + \frac{1}{d(v)}.$$

□

Since the Green values are relatively small (e.g., of order  $1/n^c$ , for some positive constant  $c$ ), we need very sharply approximate PageRank to be within a factor of  $1 + O(n^{-c})$  of the exact values in the analysis of the performance bound for the graph sparsification algorithms that we will examine in Section 4.

For all the pairs  $(u, v)$ , we define the approximate Green value  $\tilde{g}_\beta$  by

$$\tilde{g}_\beta(u, v) = \frac{\text{pr}_{\beta, \chi_u - r_{\chi_u}}(u)}{d(u)} - \frac{\text{pr}_{\beta, \chi_u - r_{\chi_u}}(v)}{d(v)} + \frac{\text{pr}_{\beta, \chi_v - r_{\chi_v}}(v)}{d(v)} - \frac{\text{pr}_{\beta, \chi_v - r_{\chi_v}}(u)}{d(u)}$$

where  $\text{pr}_{\beta, \chi_u - r_{\chi_u}}$  and  $\text{pr}_{\beta, \chi_v - r_{\chi_v}}$  are the approximate PageRank vectors as outputs of ApproximatePR for  $\text{pr}_{\beta, u}$  and  $\text{pr}_{\beta, v}$  respectively, and  $r_{\chi_u}$  and  $r_{\chi_v}$  are the corresponding residual vectors such that  $\|r_{\chi_u} D^{-1}\|_1 \leq \epsilon/4$  and  $\|r_{\chi_v} D^{-1}\|_1 \leq \epsilon/4$ . In the following, we will prove that

**Lemma 4.** *For two distinct vertices  $u, v \in V$ , we have*

$$|g_\beta(u, v) - \tilde{g}_\beta(u, v)| \leq \epsilon.$$

*Proof.* Let  $T_\beta = \beta D^{-1/2} \mathcal{G}_\beta D^{1/2}$ , and we may express  $\tilde{g}_\beta(u, v)$  as the following form

$$\begin{aligned} \tilde{g}_\beta(u, v) &= (\chi_u - r_{\chi_u}) T_\beta D^{-1} \chi_u^T - (\chi_u - r_{\chi_u}) T_\beta D^{-1} \chi_v^T + \\ &\quad (\chi_v - r_{\chi_v}) T_\beta D^{-1} \chi_v^T - (\chi_v - r_{\chi_v}) T_\beta D^{-1} \chi_u^T, \end{aligned}$$

which implies that

$$g_\beta(u, v) - \tilde{g}_\beta(u, v) = r_{\chi_u} T_\beta D^{-1} \chi_u^T + r_{\chi_v} T_\beta D^{-1} \chi_u^T + r_{\chi_v} T_\beta D^{-1} \chi_v^T + r_{\chi_u} T_\beta D^{-1} \chi_v^T.$$

Since  $\mathbf{1}Z = \mathbf{1}$ , we have  $\mathbf{1}DT_\beta D^{-1} = \mathbf{1}$  and  $\|DT_\beta D^{-1}\|_1 = 1$ . Therefore

$$\|r_{\chi_u} T_\beta D^{-1}\|_1 = \|r_{\chi_u} D^{-1} DT_\beta D^{-1}\|_1 \leq \|r_{\chi_u} D^{-1}\|_1 \|DT_\beta D^{-1}\|_1 = \epsilon/4.$$

Thus,

$$\begin{aligned} |g_\beta(u, v) - \tilde{g}_\beta(u, v)| &\leq |r_{\chi_u} T_\beta D^{-1} \chi_u^T| + |r_{\chi_v} T_\beta D^{-1} \chi_u^T| + |r_{\chi_v} T_\beta D^{-1} \chi_v^T| + |r_{\chi_u} T_\beta D^{-1} \chi_v^T| \\ &\leq \epsilon/4 + \epsilon/4 + \epsilon/4 + \epsilon/4 = \epsilon. \end{aligned}$$

The Lemma is proved. □

Here we will give a rough estimate for computing Green's values by directly using Theorem 1. Note that by invoking Theorem similarly without using further techniques the running time does not seem to improve.

**Theorem 3.** *Given any constant  $\epsilon > 0$  and any pair  $(u, v) \in V \times V$ , the approximate Green value  $\tilde{g}_\beta(u, v)$  can be computed in  $O(\frac{2+\beta}{\beta\epsilon})$  time and*

$$|g_\beta(u, v) - \tilde{g}_\beta(u, v)| \leq \epsilon.$$

*In particular, after  $O(\frac{(2+\beta)n}{\beta\epsilon})$  preprocessing time, for each  $(u, v) \in V \times V$ , we can compute such  $\tilde{g}_\beta(u, v)$  by using a constant number of queries.*

*Proof.* By Theorem 1, we can compute approximate PageRank vectors  $\text{pr}_{\beta, \chi_v - r_{\chi_v}}$  such that the residual vector  $r_{\chi_v}$  satisfies  $|r_{\chi_v}(v)/d(v)| \leq \epsilon/4$  for all  $v \in V$  in  $O(\frac{2+\beta}{\beta\epsilon})$  time. From (4), the theorem follows. □

A direct consequence of the above theorem is the following. Let  $\Delta = \max_{v \in V} d(v)$  denote the maximum degree.

**Corollary 2.** *Given any constant  $\epsilon > 0$  and any pair  $(u, v) \in V \times V$ , we can compute quantity  $\tilde{g}_\beta(u, v)$  in  $O(\frac{\Delta(2+\beta)^2}{\beta^2 \epsilon})$  time such that*

$$|g_\beta(u, v) - \tilde{g}_\beta(u, v)| \leq \epsilon g_\beta(u, v).$$

*In particular, after  $O(\frac{\Delta(2+\beta)^2 n}{\beta^2 \epsilon})$  preprocessing time, for each  $(u, v) \in V \times V$ , we can compute such  $\tilde{g}_\beta(u, v)$  by using a constant number of queries.*

*Proof.* In Lemma 3, a lower bound for  $g_\beta(u, v)$  is established:  $g_\beta(u, v) \geq \frac{\beta}{(2+\beta)\Delta}$ . Thus, for the targeted error bound  $\epsilon$ , we set  $\epsilon' = \frac{\beta}{(2+\beta)\Delta} \epsilon$  and apply Theorem 3 with parameter  $\epsilon'$  instead of  $\epsilon$ .  $\square$

We will improve both Theorem 3 and Corollary 2 in the Section 5 by using sharp approximate PageRank algorithms and dimension reduction techniques.

## 4 Graph sparsification using Green values

To construct our sparsifier, we use a method quite similar to the scheme by Spielman and Srivastava except that PageRank is used here instead of effective resistance. We will give three graph sparsification algorithms, two of which involve approximate Green values (which will be examined in section 5). In the section, we use exact Green values for edges.

Recall that the graph  $G = (V, E, w)$  we consider here is an undirected weighted graph. For a subset  $S$  of vertices in  $G$ , the *edge boundary*  $\partial(S)$  of  $S$  consists of all edges with exactly one endpoint in  $S$ . The weight of  $\partial(S)$ , denoted by  $w(S, \bar{S})$ , is the sum of all edge weights of edges in  $\partial(S)$ . The *volume* of  $S$ , denoted by  $\text{vol}(S)$ , is defined to be the sum of *degrees*  $d(v)$  over all  $v$  in  $S$ . When  $S = V$ , we write  $\text{vol}(S) = \text{vol}(G)$ . The *Cheeger ratio* (or *conductance*)  $h_G(S)$  of  $S$  in  $G$  is defined by

$$h_G(S) = \frac{w(S, \bar{S})}{\min\{\text{vol}(S), \text{vol}(\bar{S})\}}.$$

The *conductance*  $h_G$  of  $G$  is defined to be the minimum Cheeger ratio among all subsets  $S$  with  $\text{vol}(S) \leq \text{vol}(G)$ .

The goal of sparsification is to approximate a given graph  $G$  by a sparse graph  $\tilde{G}$  on the same set of vertices while the sparse graph  $\tilde{G}$  preserves the Cheeger ratios of every subset of vertices to within a factor of  $1 + \epsilon$ . The sparse graph  $\tilde{G}$  is called a *sparsifier* of  $G$  if the Laplacian  $L$  and  $\tilde{L}$  for  $G$  and  $\tilde{G}$  satisfy (1) for all  $x \in \{0, 1\}^n$ .

The main step in any sparsification algorithm [9, 30–32, 46–48] is to choose an appropriate probability distribution for random sampling the edges in a way that Cheeger ratios of subsets change little. Our sparsification algorithm is a sampling process using probabilities proportional to the Green values  $g_\beta$ 's as follows:

$\tilde{G} = \text{SparsifyExactGreen}(G, q, \beta)$ :

For each  $e = (u, v) \in E$ , set probability  $p_e \propto w(e)g_\beta(e)$  and repeat the following steps for  $q$  times:

1. Choose an edge  $e \in G$  randomly with probability  $p_e$
2. Add  $e$  to  $\tilde{G}$  with weight  $w(e)/qp_e$ .
3. Sum the weights if an edge is chosen more than once.

The analysis of the above algorithm will be examined in the following subsections. Our main theorem is the following

**Theorem 4.** *In an unweighted graph  $G$  with  $m$  edges, for any  $1 > \epsilon > 0$ , let  $\tilde{G}$  denote the output of the algorithm  $\text{SparsifyExactGreen}(G, q, \beta)$ , where  $q = 256C^2n \log n / \epsilon^2$ ,  $\beta = 1/2$ ,  $C \geq 1$  is a absolute constant. Then with probability at least  $1/2$ , we have*

1. (Performance guarantee) For all  $S \in V$ ,  $|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon$ .
2. (Sampling complexity) Algorithm  $\text{SparsifyExactGreen}$  can be performed by using  $O(\frac{n \log n}{\epsilon^2})$  sampling.

#### 4.1 Analyzing the sparsifier

Our analysis follows the general scheme as that of [48]. In our analysis of sparsifier, we consider the matrix  $\Lambda_\beta = W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2}$ . Note that  $\Lambda_\beta$  is a  $(n+m) \times (n+m)$  matrix and we index its the first  $n$  columns (rows) by  $V$  and its last  $m$  columns (rows) by  $E$ . From the definition of Green values in (3) and (4), it is easy to verify that  $\Lambda_\beta(e, e) = \frac{1}{\beta} \sqrt{W_\beta(e, e)} g_\beta(e) \sqrt{W_\beta(e, e)} = \frac{1}{\beta} w(e) g_\beta(e)$ . Here are several useful properties for  $\Lambda_\beta$ .

- Lemma 5.** (i)  $\Lambda_\beta^2 = \Lambda_\beta$ .  
(ii) The dimension (or rank) of  $\Lambda_\beta$ , denoted by  $\dim(\Lambda_\beta)$  is  $n$ .  
(iii) The eigenvalues of  $\Lambda_\beta$  are 1 with multiplicity  $n$  and 0 with multiplicity  $m$ .  
(iv)  $\Lambda_\beta(e, e) = \|\Lambda_\beta(\cdot, e)\|_2^2$ .

*Proof.* For (i), we note that

$$\begin{aligned} \Lambda_\beta \Lambda_\beta &= (W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2}) (W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2}) \\ &= W_\beta^{1/2} S' \mathcal{G}_\beta (S'^T W_\beta S') \mathcal{G}_\beta S'^T W_\beta^{1/2} \\ &= W_\beta^{1/2} S' \mathcal{G}_\beta \mathcal{L}_\beta \mathcal{G}_\beta S'^T W_\beta^{1/2} \\ &= W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2} \\ &= \Lambda_\beta. \end{aligned}$$

For (ii), it is easy to verify that  $\dim(W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2}) = \dim(S' \mathcal{G}_\beta S'^T) = \dim(\mathcal{G}_\beta) = n$ .

For (iii), Since  $\Lambda_\beta^2 = \Lambda_\beta$ , the eigenvalue of  $\Lambda_\beta$  are all 0 or 1. Since  $\dim(\Lambda_\beta) = n$ , there must be  $n$  nonzero eigenvalues.

(iv) follows from  $\Lambda_\beta$  is symmetric. □

Next, we introduce some notations and several lemmas that the theorems in later sections rely on. Let  $\tilde{w}(e)$  be the edge weight of edge  $e$  in  $\tilde{G}$ . Recall that  $q$  is a number and  $p_e$  is the sampling probability for  $e \in E$ . Denote  $I_\beta$  as a nonnegative diagonal matrix

$$I_\beta = \begin{bmatrix} I_{n \times n} & 0 \\ 0 & R \end{bmatrix}_{(n+m) \times (n+m)},$$

where  $R$  is a  $m \times m$  nonnegative diagonal matrix denoted by

$$R(e, e) = \frac{\tilde{w}(e)}{w(e)} = \frac{\# \text{ of times } e \text{ is sampled}}{qp_e}.$$

**Lemma 6.** *Suppose  $I_\beta$  is a nonnegative diagonal matrix such that  $\|\Lambda_\beta I_\beta \Lambda_\beta - \Lambda_\beta \Lambda_\beta\|_2 \leq \epsilon$ . Then*

$$\forall x \in \mathbb{R}^n, \quad |x \tilde{\mathcal{L}}_\beta x^T - x \mathcal{L}_\beta x^T| \leq \epsilon x \mathcal{L}_\beta x^T,$$

where  $\mathcal{L}_\beta = S'^T W_\beta S'$  and  $\tilde{\mathcal{L}}_\beta = S'^T W_\beta^{1/2} I_\beta W_\beta^{1/2} S'$ .

*Proof.* The assumption of this lemma is equivalent to

$$\sup_{y \in \mathbb{R}^{m+n}, y \neq 0} \frac{|y \Lambda_\beta (I_\beta - I) \Lambda_\beta y^T|}{yy^T} \leq \epsilon,$$

which implies

$$\sup_{y^T \in \text{im}(W_\beta^{1/2} S'), y \neq 0} \frac{|y \Lambda_\beta (I_\beta - I) \Lambda_\beta y^T|}{yy^T} \leq \epsilon.$$

Here  $y \in \text{im}(M)$  means  $y = xM$  for some  $x$ . We have

$$\begin{aligned} & \sup_{y^T \in \text{im}(W_\beta^{1/2} S'), y \neq 0} \frac{|y \Lambda_\beta (I_\beta - I) \Lambda_\beta y^T|}{yy^T} \\ = & \sup_{x \in \mathbb{R}^n, W_\beta^{1/2} S' x^T \neq 0} \frac{|x S'^T W_\beta S' \mathcal{G}_\beta S'^T W_\beta^{1/2} (I_\beta - I) W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta S' x^T|}{x S'^T W_\beta S' x^T} \\ = & \sup_{x \in \mathbb{R}^n, W_\beta^{1/2} S' x^T \neq 0} \frac{|x \tilde{\mathcal{L}}_\beta \mathcal{G}_\beta S'^T W_\beta^{1/2} (I_\beta - I) W_\beta^{1/2} S' \mathcal{G}_\beta \mathcal{L}_\beta x^T|}{x S'^T W_\beta S' x^T} \\ = & \sup_{x \in \mathbb{R}^n, W_\beta^{1/2} S' x^T \neq 0} \frac{|x S'^T W_\beta^{1/2} (I_\beta - I) W_\beta^{1/2} S' x^T|}{x S'^T W_\beta S' x^T} \\ = & \sup_{x \in \mathbb{R}^n, W_\beta^{1/2} S' x^T \neq 0} \frac{|x \tilde{\mathcal{L}}_\beta x^T - x \mathcal{L}_\beta x^T|}{x \mathcal{L}_\beta x^T} \leq \epsilon \end{aligned}$$

The lemma follows from the fact that  $\dim(\text{im}(W_\beta^{1/2} S')) = n$  and  $W_\beta^{1/2} S' x^T = 0$  if and only if  $x = 0$ .  $\square$

**Lemma 7 ([44]).** *Let  $\mathbf{p}$  be a probability distribution over  $\Omega \subseteq \mathbb{R}^d$  such that  $\sup_{y \in \Omega} \|y\|_2 \leq M$  and  $\mathbb{E}_{\mathbf{p}} \|y^T y\|_2 \leq 1$ . Let  $y_1 \dots y_q$  be independently samples drawn from  $\mathbf{p}$ . Then for every  $1 > \epsilon > 0$ ,*

$$\mathbb{P} \left\{ \left\| \frac{1}{q} \sum_{i=1}^q y_i^T y_i - \mathbb{E} y^T y \right\|_2 > \epsilon \right\} \leq 2 \exp(-\epsilon^2/a^2),$$

where  $a = \min \left( CM \sqrt{\frac{\log q}{q}}, 1 \right)$  and  $C$  is an absolute constant.

## 4.2 Proof of Theorem 4

We first prove the following theorem which leads to the proof of Theorem 4.

**Theorem 5.** *Let  $\mathcal{L}$  be the normalized Laplacian of  $G$ . Let  $\tilde{G}$  be the output of the algorithm **SparsifyExactGreen**( $G, q, \beta$ ), where  $q = 4C^2 n \log n / \epsilon^2$ ,  $1 \geq \epsilon > 0$  and  $C$  is a absolute constant. Then with probability at least  $1/2$ , we have*

$$\forall x \in \mathbb{R}^n, \quad |x \tilde{\mathcal{L}}_\beta x^T - x \mathcal{L}_\beta x^T| \leq \epsilon x \mathcal{L}_\beta x^T, \quad (5)$$

where  $\mathcal{L}_\beta = \beta I + \mathcal{L} = S'^T W_\beta S'$  and  $\tilde{\mathcal{L}}_\beta = S'^T W_\beta^{1/2} I_\beta W_\beta^{1/2} S'$ .

*Remark 1.* Let  $\mathcal{L}'$  be the matrix  $D^{-1/2} \tilde{D}^{1/2} \tilde{\mathcal{L}} \tilde{D}^{1/2} D^{-1/2}$  where  $\tilde{D}$  and  $\tilde{\mathcal{L}}$  are the diagonal matrices of weighted degree and the normalized Laplacian of graph  $\tilde{G}$ , respectively. We observe that the inequality in (5) is equivalent to the following:

$$\forall x \in \mathbb{R}^n, \quad |x \mathcal{L}' x^T - x \mathcal{L} x^T| \leq \epsilon x \mathcal{L}_\beta x^T. \quad (6)$$

*Remark 2.* The probability  $1/2$  can be improved to “high probability”, i.e.,  $1 - 1/n^{O(1)}$ . However, the number of sampled edges will be increased from  $O(n \log n)$  to  $O(n \log^2 n)$ .

*Proof of Theorem 5:* Before applying Lemma 7, we observe that

$$\begin{aligned}\Lambda_\beta I_\beta \Lambda_\beta &= \sum_{e \in E} R(e, e) \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) + \sum_{v \in V} \Lambda_\beta(v, \cdot)^T \Lambda_\beta(v, \cdot) \\ \Lambda_\beta \Lambda_\beta &= \sum_{e \in E} \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) + \sum_{v \in V} \Lambda_\beta(v, \cdot)^T \Lambda_\beta(v, \cdot).\end{aligned}$$

Thus we have

$$\Lambda_\beta I_\beta \Lambda_\beta - \Lambda_\beta \Lambda_\beta = \sum_{e \in E} (R(e, e) - 1) \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot).$$

Now, we consider

$$\begin{aligned}& \sum_{e \in E} R(e, e) \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) \\ &= \sum_{e \in E} \frac{\# \text{ of times } e \text{ is sampled}}{qp_e} \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) \\ &= \frac{1}{q} \sum_{e \in E} (\# \text{ of times } e \text{ is sampled}) \frac{\Lambda_\beta(e, \cdot)^T}{\sqrt{p_e}} \frac{\Lambda_\beta(e, \cdot)}{\sqrt{p_e}} \\ &= \frac{1}{q} \sum_{i=1}^q y_i^T y_i\end{aligned}$$

where  $y_1, \dots, y_q$  are random vectors drawn independently with replacement from the distribution  $\mathbf{p}$  defined by setting  $y = \frac{1}{\sqrt{p_e}} \Lambda_\beta(e, \cdot)$  with probability  $p_e$ .

We also need to bound the norm of the expectation of  $y^T y$  as follows:

*Claim.*  $\|\mathbb{E}_{\mathbf{p}} y^T y\|_2 \leq 1$ .

*Proof.*

$$\begin{aligned}\|\mathbb{E}_{\mathbf{p}} y^T y\|_2 &= \left\| \sum_{e \in E} p_e \frac{1}{p_e} \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) \right\|_2 \\ &= \left\| \sum_{e \in E} \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) \right\|_2 \\ &= \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{\sum_{e \in E} x \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) x^T}{x x^T}\end{aligned}$$

Since for any  $x \in \mathbb{R}^{m+n}$ , we have  $\sum_{v \in V} x \Lambda_\beta(v, \cdot)^T \Lambda_\beta(v, \cdot) x^T \geq 0$ . Thus

$$\begin{aligned}
\|\mathbb{E}_{\mathbf{p}} y^T y\|_2 &= \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{\sum_{e \in E} x \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) x^T}{x x^T} \\
&\leq \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{\sum_{e \in E} x \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) x^T + \sum_{v \in V} x \Lambda_\beta(v, \cdot)^T \Lambda_\beta(v, \cdot) x^T}{x x^T} \\
&= \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{x (\sum_{e \in E} \Lambda_\beta(e, \cdot)^T \Lambda_\beta(e, \cdot) + \sum_{v \in V} \Lambda_\beta(v, \cdot)^T \Lambda_\beta(v, \cdot)) x^T}{x x^T} \\
&= \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{x \Lambda_\beta \Lambda_\beta x^T}{x x^T} \\
&= \sup_{x \in \mathbb{R}^{m+n}, x \neq 0} \frac{x \Lambda_\beta x^T}{x x^T} \\
&= \|\Lambda_\beta\|_2 = 1.
\end{aligned}$$

The last equality above follows from the fact that the eigenvalues of  $\Lambda_\beta$  are all 1 or 0, and the claim follows.

An upper bound for the norm of  $y$  can be established as follows:

$$\frac{1}{\sqrt{p_e}} \|\Lambda_\beta(e, \cdot)\|_2 = \frac{1}{\sqrt{p_e}} \sqrt{\Lambda_\beta(e, e)} = \sqrt{\frac{\sum_{e' \in E} w(e') g_\beta(e')}{w_e g_\beta(e)}} \sqrt{\frac{w_e g_\beta(e)}{\beta}} = \sqrt{\sum_{e' \in E} \frac{w(e') g_\beta(e')}{\beta}}.$$

Note that  $\Lambda_\beta = W_\beta^{1/2} S' \mathcal{G}_\beta S'^T W_\beta^{1/2}$  can be expressed as the following form.

$$\Lambda_\beta \triangleq \begin{bmatrix} \beta \mathcal{G}_\beta & \beta^{1/2} \mathcal{G}_\beta S'^T W_\beta^{1/2} \\ \beta^{1/2} W_\beta^{1/2} S \mathcal{G}_\beta & W_\beta^{1/2} S \mathcal{G}_\beta S'^T W_\beta^{1/2} \end{bmatrix}_{(n+m) \times (n+m)}.$$

Thus,

$$\sum_{e \in E} \frac{w_e g_\beta(e)}{\beta} = \sum_{e \in E} \Lambda_\beta(e, e) = \text{Tr}(W_\beta^{1/2} S \mathcal{G}_\beta S'^T W_\beta^{1/2}) = \text{Tr}(\Lambda_\beta) - \text{Tr}(\beta \mathcal{G}_\beta).$$

Since  $\text{Tr}(\Lambda_\beta) = n$ , we have  $\frac{1}{\sqrt{p_e}} \|\Lambda_\beta(e, \cdot)\|_2 \leq \sqrt{n}$ . By setting  $q = 4C^2 n \log n / \epsilon^2$ , we have

$$\min \left( CM \sqrt{\frac{\log q}{q}}, 1 \right) \leq C \sqrt{\epsilon^2 \frac{n \log(4C^2 n \log n / \epsilon^2)}{4C^2 n \log n}} \leq \epsilon/2. \quad (7)$$

By applying the Rudelson and Vershynin's lemma in [44] (Lemma 7), we completes the proof of the theorem.  $\square$

Before applying Theorem 5 to prove Theorem 4, we still need the following two lemmas. We will here consider  $G$  as an unweighted graph first, i.e  $w(e) = 1$  for all edges, although this can be easily extended to the general weighted graphs.

**Lemma 8.** *Let  $\tilde{G}$  be the output of algorithm *SparsifyExactGreen*( $G, q, \beta$ ), where  $q = 4C^2 n(\beta + 2) \log n / \epsilon^2$ . Then, with probability  $1 - 1/n$ , for all subsets  $S \subset V$ , we have  $|\text{vol}_{\tilde{G}}(S) - \text{vol}_G(S)| \leq \epsilon \text{vol}_G(S)$ .*

*Proof.* For a fixed set  $S \subset V$ , let  $E(S)$  be the set of edges with at least one endpoint in  $S$ . Consider the i.i.d. random variables  $X_1, X_2, \dots, X_q$  defined as follows:

$$X_i = \begin{cases} 1/p_e & \text{with probability } p_e \text{ for } e \in E(S) \\ 0 & \text{otherwise.} \end{cases}$$

We consider random variable  $V_S = \frac{1}{q} \sum_{i=1}^q X_i$ . It is easy to verify that

$$V_S = \sum_{e \in E(S)} \frac{\# \text{ of times } e \text{ is sampled}}{qp_e} = \sum_{e \in E(S)} \frac{\tilde{w}(e)}{w(e)} = \text{vol}_{\tilde{G}}(S).$$

Clearly,

$$\mathbb{E}[X_i] = \sum_{e \in E(S)} p_e \frac{1}{p_e} = \text{vol}_G(S),$$

and

$$\mathbb{E}[X_i^2] = \sum_{e \in E(S)} p_e \frac{1}{p_e^2} = \sum_{e \in E(S)} \frac{1}{p_e} = \sum_{e \in E(S)} \frac{\sum_{e' \in E} w(e') g_\beta(e')}{w(e) g_\beta(e)}.$$

Note that  $\sum_{e \in E} w_e g_\beta(e) \leq \beta \text{Tr}(A_\beta) \leq \beta n$ . Using Lemma 3, we have

$$\begin{aligned} \mathbb{E}[X_i^2] &= \sum_{e \in E(S)} \frac{\sum_{e' \in E} w(e') g_\beta(e')}{w(e) g_\beta(e)} \\ &\leq \sum_{e \in E(S)} \frac{\beta n}{g_\beta(e)} \\ &\leq \sum_{e=(u,v) \in E(S)} \frac{n\beta(\beta+2) \min\{d(u), d(v)\}}{\beta} \\ &\leq n(\beta+2) \text{vol}^2(S). \end{aligned}$$

By the Bernstein inequality, for i.i.d. random variables  $X_1, X_2, \dots, X_k$  such that  $|X_i| \leq M$ .

$$\mathbb{P} \left\{ \left| \frac{1}{k} \sum_{i=1}^k X_i - \frac{1}{k} \sum_{i=1}^k \mathbb{E}[X_i] \right| > \frac{t}{k} \sum_{i=1}^k \mathbb{E}[X_i] \right\} \leq \exp \left\{ - \frac{t^2 (\sum_{i=1}^k \mathbb{E}[X_i])^2}{\sum_{i=1}^k \mathbb{E}[X_i^2] + Mt/3} \right\}.$$

In our case,  $t = \epsilon$ ,  $k = q$ ,  $\sum_{i=1}^q \mathbb{E}[X_i] = q \text{vol}_G(S)$ ,  $C > 1$ , and  $M = n(\beta+2)\Delta_S$  where  $\Delta_S$  is the maximum degree of vertices in  $S$ , i.e.  $\Delta_S = \max_{v \in S} d(v)$ . Thus,

$$\begin{aligned} \exp \left\{ - \frac{t^2 (\sum_{i=1}^k \mathbb{E}[X_i])^2}{\sum_{i=1}^k \mathbb{E}[X_i^2] + Mt/3} \right\} &\leq \exp \left\{ - \frac{q^2 \epsilon^2 \text{vol}^2(S)}{qn(\beta+2) \text{vol}^2(S) + n\epsilon(\beta+2)\Delta_S/3} \right\} \\ &\leq \exp \left\{ - \frac{qt^2}{2(\beta+2)n} \right\} \\ &= \exp \left\{ - \frac{2\epsilon^2 C^2 \log n}{\epsilon^2} \right\} \\ &\leq 1/n^2. \end{aligned} \tag{8}$$

Therefore, the probability of the event that there exists a vertex  $v$  such that  $|\text{vol}_{\tilde{G}}(v) - \text{vol}_G(v)| > \epsilon$  is less than  $1/n$  and

$$|\text{vol}_{\tilde{G}}(S) - \text{vol}_G(S)| \leq \sum_{v \in V} |\text{vol}_{\tilde{G}}(v) - \text{vol}_G(v)| \leq \epsilon \sum_{v \in V} \text{vol}_G(v) \leq \epsilon \text{vol}_G(S).$$

The proof of the lemma is complete □

**Lemma 9.** *If sparse graph  $\tilde{G}$  corresponding to graph  $G$  satisfies the following two conditions:*

1. for all  $x \in \mathbb{R}^n$ ,  $|x\tilde{\mathcal{L}}_\beta x - x\mathcal{L}_\beta x^T| \leq \epsilon x\mathcal{L}_\beta x^T$ ;
2. for all subsets  $S \subset V$ ,  $|\text{vol}_{\tilde{G}}(S) - \text{vol}_G(S)| \leq \epsilon \text{vol}_G(S)$ .

Then  $|h_{\tilde{G}}(S) - h_G(S)| \leq 2\epsilon h_G(S) + \epsilon\beta$ .

*Proof.* Let  $L$  and  $\tilde{L}$  denote the Laplacians of the graphs  $G$  and  $\tilde{G}$  respectively. Let  $D$  and  $\tilde{D}$  be the weighted degree matrices of the graphs  $G$  and  $\tilde{G}$  respectively. For  $\beta > 0$ , we consider  $L_\beta = \beta D + L$  and  $\tilde{L}_\beta = \beta \tilde{D} + \tilde{L}$ . Let  $\chi_S$  denote the characteristic function on  $S \in V$ , satisfying  $\chi_S(v) = 1$  if  $v \in S$  and  $\chi_S(v) = 0$  otherwise. From the definition of Cheeger ratio, we have

$$\begin{aligned} |h_{\tilde{G}}(S) - h_G(S)| &= \left| \frac{\chi_S \tilde{L}_\beta \chi_S^T}{\chi_S D \chi_S^T} - \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \\ &\leq \left| \frac{\chi_S \tilde{L}_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \left| \frac{\chi_S (D - \tilde{D}) \chi_S^T}{\chi_S \tilde{D} \chi_S^T} \right| + \left| \frac{\chi_S \tilde{L}_\beta \chi_S^T}{\chi_S D \chi_S^T} - \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right|. \end{aligned}$$

Recall that  $\mathcal{L}_\beta = S'^T W_\beta S'$  and  $\tilde{\mathcal{L}}_\beta = S'^T W_\beta^{1/2} I_\beta W_\beta^{1/2} S'$ . It is straightforward to verify that

$$\frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} = \frac{\chi_S D^{1/2} \mathcal{L}_\beta D^{1/2} \chi_S^T}{\chi_S D \chi_S^T} \quad \text{and} \quad \frac{\chi_S \tilde{L}_\beta \chi_S^T}{\chi_S D \chi_S^T} = \frac{\chi_S D^{1/2} \tilde{\mathcal{L}}_\beta D^{1/2} \chi_S^T}{\chi_S D \chi_S^T}.$$

Thus, the first condition implies that

$$\left| \frac{\chi_S \tilde{L}_\beta \chi_S^T}{\chi_S D \chi_S^T} - \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \leq \epsilon \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right|.$$

Now, by the second condition, we have

$$\begin{aligned} |h_{\tilde{G}}(S) - h_G(S)| &\leq (1 + \epsilon) \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \left| \frac{\chi_S (D - \tilde{D}) \chi_S^T}{\chi_S \tilde{D} \chi_S^T} \right| + \epsilon \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \\ &\leq (1 + \epsilon) \epsilon \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| + \epsilon \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \\ &\leq 2\epsilon \left| \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} \right| \\ &\leq 2\epsilon \frac{\chi_S L_\beta \chi_S^T}{\chi_S D \chi_S^T} + \epsilon\beta \\ &= 2\epsilon h_G(S) + \epsilon\beta. \end{aligned}$$

The lemma follows.  $\square$

*Proof of Theorem 4:* To prove Theorem 4, we combine Lemma 8, Lemma 9 and Theorem 5. For any  $1 > \epsilon > 0$ , let  $\tilde{G}$  be the output of the algorithm **SparsifyExactGreen**( $G, q, \beta$ ), where  $q = 256C^2 n \log n / \epsilon^2$ ,  $\beta = 1/2$ , and  $C \geq 1$  is a constant. By Theorem 5 and Lemma 8, the conditions of Lemma 9 are satisfied with probability at least  $1/2$ . Note that we have chosen  $\beta$  to be  $1/2$  and  $h_G(S) \leq 1$ , thus algorithm **SparsifyExactGreen** can be applied by using  $O(\frac{n \log n}{\epsilon^2})$  sampling. Furthermore, for all  $S \in V$ , we have

$$|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon.$$

$\square$

By choosing a different  $\beta$ , namely,  $\beta = \phi/2$ , we have the following:

**Theorem 6.** *Given any  $1 \geq \epsilon > 0$  and  $1 \geq \phi > 0$ , let  $\tilde{G}$  be the output of the algorithm **SparsifyExactGreen**( $G, q, \beta$ ), where  $q = 256C^2 n \log n / \epsilon^2$ ,  $\beta = \phi/2$ ,  $C \geq 1$ , and  $G$  is an unweighted graph with  $m$  edges. Then with probability at least  $1/2$ , we have*

1. (Performance guarantee) For all  $S \in V$  with  $h_G(S) \leq \phi$ ,  $|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon h_G(S)$ .
2. (Sampling complexity) Algorithm **SparsifyExactGreen** can be performed by using  $O(\frac{n \log n}{\epsilon^2})$  sampling.

## 5 Sparsification using approximate PageRank vectors

In Corollary 2, we compute approximate Green values  $\tilde{g}_\beta(u, v)$  for all edges  $(u, v) \in E$  satisfying

$$(1 - \kappa)g_\beta(u, v) \leq \tilde{g}_\beta(u, v) \leq (1 + \kappa)g_\beta(u, v)$$

in  $O(\Delta n)$  time, where  $\kappa$  is any absolute constant such that  $\kappa < 1$  (e.g.,  $\kappa = 0.01$ ). Instead of exact Green values, we can use approximate Green values as we run the algorithm **SparsifyExactGreen**. The approximate Green values  $\tilde{g}_\beta$ 's are combinations of approximate PageRank vectors  $\text{pr}'_{\beta, v}$ 's. Here we choose the following parameters for **ApproximatePR** :

$$\text{pr}'_{\beta, v} = \mathbf{ApproximatePR}(\chi_v, \beta, \frac{\beta}{(2 + \beta)\Delta}\kappa).$$

It is not difficult to verify that all results in Section 4 will change at most by a constant factor if we run the algorithm **SparsifyExactGreen** by using approximate Green values as above. The performance guarantee and the number of sampled edges in the Theorems 4 differ by at most a constant factor, although the computational complexity will increase to  $O(\Delta n)$ . In order to further improve the running time, we use several methods in the following subsections.

### 5.1 Graph sparsification by using sharply approximate Green values

In order to have better error estimate of approximate Green values, we need to improve the error estimate for approximate PageRank vectors in Theorem 1. We will use the strengthened approximate PageRank algorithm **SharpApproximatePR** and the dimension reduction technique in [48] to approximate the Green values by using these sharply approximate PageRank vectors produced by **SharpApproximatePR**.

First, we notice that

$$\begin{aligned} g_\beta(u, v) &= \beta(\chi_u - \chi_v)D^{-1/2}\mathcal{G}_\beta D^{-1/2}(\chi_u - \chi_v)^T \\ &= \beta(\chi_u - \chi_v)D^{-1/2}\mathcal{G}_\beta \mathcal{L}_\beta \mathcal{G}_\beta D^{-1/2}(\chi_u - \chi_v)^T \\ &= \beta \|W_\beta S' \mathcal{G}_\beta D^{-1/2}(\chi_u - \chi_v)^T\|_2^2 \\ &= \frac{1}{\beta} \|W_\beta S' D^{1/2}[\beta D^{-1/2} \mathcal{G}_\beta D^{-1/2}](\chi_u - \chi_v)^T\|_2^2. \end{aligned}$$

Therefore,  $g_\beta(u, v)$ 's are just pairwise distances between vectors  $\{\mathcal{Z}\chi_v^T\}_{v \in V}$  where

$$\mathcal{Z} = W_\beta S' D^{1/2}[\beta D^{-1/2} \mathcal{G}_\beta D^{-1/2}].$$

However, the dimension of the vectors in  $\{\mathcal{Z}\chi_v^T\}_{v \in V}$  is  $m+n$ . In order to reduce the computational complexity for computing these vectors, we project these vectors into a lower dimensional space while preserving their pairwise distances by the following lemma.

**Lemma 10 ([1]).** *Given vectors  $x_1, \dots, x_n \in \mathbb{R}^d$  and constants  $\epsilon, \gamma > 0$ , let  $k_0 = c_\gamma \log n / \epsilon^2$  where  $c_\gamma$  is a constant depending on  $\gamma$ . For integer  $k \geq k_0$ , let  $R_{k \times d}$  be a random matrix where  $\{R_{ij}\}$  are independent random variables with values  $\pm 1/\sqrt{k}$ . Then with probability  $1 - \frac{1}{n^\gamma}$ , we have*

$$(1 - \epsilon)\|x_i - x_j\|_2^2 \leq \|Rx_i - Rx_j\|_2^2 \leq (1 + \epsilon)\|x_i - x_j\|_2^2.$$

Now, we are ready to state our algorithm to approximate the Green values.

**ApproxiamteGreen**( $\beta, \epsilon, k$ ):

1. Let  $R_{k \times (n+m)} = [R_1, R_2]$  be a random matrix whose entries are independent random variables with values  $\pm 1/\sqrt{k}$ , where  $R_1$  is an  $k \times n$  matrix and  $R_2$  is an  $k \times m$  matrix.
2. Let  $Y = RW_\beta^{1/2} S' D^{1/2}$  and  $\tilde{Z} = RZ$ .
3. For  $i = 1, \dots, k$ , do the following
  - (a) Let  $y_i$  be the  $i$ th row of  $Y$  and  $\tilde{z}_i$  be the  $i$ th row of  $\tilde{Z}$ .
  - (b) Approximate  $\tilde{z}_i$  by calling  $\tilde{z}_i' = \mathbf{SharpApproximatePR}(y_i, \beta, \epsilon/n^r)$  where  $r$  is a constant.
4. Let  $\tilde{Z}'$  be the approximated matrix for  $\tilde{Z}$  whose rows are  $\tilde{z}_1', \tilde{z}_2', \dots, \tilde{z}_k'$ .  
For all  $(u, v) \in E$ , return  $\tilde{g}_\beta(u, v) = \|\tilde{Z}'(\chi_u - \chi_v)^T\|_2^2$ .

In order to analysis algorithm **ApproxiamteGreen**, we need to give a bound for  $y_i$ 's by the following lemma.

**Lemma 11.** *Given an integer  $k$  and random matrix whose entries are independent random variables with values  $\pm 1/\sqrt{k}$ , with probability  $1 - 1/n^2$ , we have  $\|y_i\|_1 \leq c(\sum_{v \in V} d(v)^{1/2})\sqrt{\log n}$  for  $1 \leq i \leq k$ , where  $c$  is an absolute constant.*

*Proof.* To simplify the argument, we consider a graph  $G$  with edge weight 1 for all edges in  $G$ . Denote the  $v$ th element of vector  $y_i$  to be  $y_{i,v}$ . From the definition of  $Y$ , we have  $Y = \beta^{1/2} R_1 D^{1/2} + R_2 W^{1/2} B$ . We now fix a vertex  $v \in V$ . Let  $X_0, X_1, \dots, X_l$  be independent random variables with values  $\pm 1/\sqrt{k}$  (i.e., with probability  $1/2$  to be  $1/\sqrt{k}$  and with probability  $1/2$  to be  $-1/\sqrt{k}$ ). We can write  $y_{i,v}$  as follows

$$y_{i,v} = \sum_{l=1}^{d(v)} X_l + \beta^{1/2} \sqrt{d(v)} X_0.$$

By Chernoff's inequality, we have

$$\mathbb{P}\left\{\left|\sum_{l=1}^{d(v)} X_l\right| > \sqrt{ka}\right\} < 2 \exp\left\{-\frac{ka^2}{2d(v)}\right\},$$

which is equivalent to

$$\mathbb{P}\left\{\sum_{l=1}^{d(v)} X_l > c\sqrt{d(v) \log n}\right\} < 2 \exp\left\{-\frac{c^2 \log n}{2}\right\}$$

where  $c \geq 4$  is a constant. Notice that  $\beta \leq 1$  and  $|X_0| = 1/\sqrt{k}$ . The lemma follows by applying the above inequality on the following event.

$$\bigcup_{v \in V} \left\{ |y_{i,v}| > c\sqrt{d(v) \log n} \right\}.$$

□

By combining the above lemmas and Theorem 2, we have the following lemma.

**Theorem 7.** *Given any constant  $\epsilon > 0$ , we set  $k = c \log n / \epsilon^2$ . If  $\beta = \Omega(\text{poly}(1/n))$ , then algorithm **ApproxiamteGreen**( $\beta, \epsilon, k$ ) returns  $\tilde{g}_\beta(u, v)$  satisfying  $|g_\beta(u, v) - \tilde{g}_\beta(u, v)| \leq \epsilon g_{\beta(u, v)}$  in  $O(\frac{m \log^2 n}{\beta \epsilon^2})$  time.*

*Proof.* To bound the running time, note that step 1 can be completed in  $O(m \log n / \epsilon)$  time since  $R$  is a  $k \times (n + m)$  random matrix. In step 2, we set  $Y = RW_\beta^{1/2} S' D^{1/2}$  and it only takes  $O(m \log n / \epsilon^2)$  time since  $S'$  has  $O(m)$  entries and  $W_\beta^{1/2}$  is a diagonal matrix.

In step 3, Let  $y_i$  be the  $i$ th row of  $Y$  and  $\tilde{z}_i$  be the  $i$ th row of  $\tilde{Z}$ , i.e.,

$$(Y[\beta D^{-1/2} \mathcal{G}_\beta D^{-1/2}])_{k \times n}.$$

Therefore, we have  $\tilde{z}_i = y_i[\beta D^{-1/2} \mathcal{G}_\beta D^{-1/2}]$  and we can view  $\tilde{z}_i$  as a scaled PageRank vector with seed vector  $y_i$ .

In Lemma 11, we have proved that with probability at least  $1 - 1/n$ ,  $\|y_i\|_1 = O(\sum_{v \in V} \sqrt{d(v) \log n})$  for  $1 \leq i \leq k$ . Without loss of generality, we may assume  $O(\sum_{v \in V} \sqrt{d(v) \log n}) = O(m)$  otherwise the graph is sufficient sparse. Thus,  $\tilde{z}_i$  can be approximated by using algorithm **ApproximatePRB** with arbitrary small absolute error, say,  $\epsilon'$ . Thus, each call of **ApproximatePRB** just takes  $O(\frac{m \log(1/\epsilon')}{\beta})$  time. By Lemma 3,  $g_\beta(u, v) = \Omega(\beta/n)$  which implies that we only need to set  $\epsilon' = \epsilon/n^r$  for some large enough but fixed constant  $r$ . Thus, each call of **ApproximatePRB** will actually take  $O(m \log n/\beta)$  time. Since there are  $k = O(\log n/\epsilon^2)$  calls, the total running time of step 3 is  $O(m \log^2 n/\beta \epsilon^2)$ .

In step 4, since each column of  $\tilde{Z}'$  has  $k = O(\log n/\epsilon^2)$  entries and there are  $m$  edges, the running time of step 4 is  $O(m \log n/\epsilon^2)$ . The lemma then follows.  $\square$

$\tilde{G} = \mathbf{SparsifyApprGreen}(G, q, \beta, \epsilon)$ :

1. For all  $(u, v) \in E$ , compute approximate Green values by **ApproximateGreen** $(\beta, \kappa, k)$  where  $\kappa = 1/2$  and  $k = c \log n/\epsilon^2$  for some absolute constant  $c$ .
2. Apply **SparsifyApprGreen** $(G, q, \beta)$  with approximate Green values  $\tilde{g}_\beta(e)$ 's.

**Theorem 8.** *Given any constant  $\epsilon > 0$ ,  $\phi > 0$  and an unweighted graph  $G$  with  $m$  edges, set  $q = 256C^2 n \log n/\epsilon^2$ ,  $\beta = \phi/2$ , and let  $\tilde{G}$  be the output of the algorithm **SparsifyApprGreen** $(G, q, \beta, \epsilon)$  where we use algorithm **ApproximateGreen** as a subroutine to approximate the Green value  $\tilde{g}_\beta(u, v)$  for all  $(u, v) \in E$ . Then with probability at least  $1/2$ , we have*

1. (Performance guarantee) For all  $S \in V$  satisfying  $h_{\tilde{G}}(S) \geq \phi$ , we have  $|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon h_G(S)$ .
2. (Complexity guarantee) Algorithm **SparsifyApprGreen** can be performed by using  $O(\frac{m \log^2 n}{\phi})$  preprocessing time, and  $O(\frac{n \log n}{\epsilon^2})$  sampling.

The proof of Theorem 8 is quite similar to the analysis in Section 4 and will be omitted.

## 6 Graph partitioning using approximate PageRank vectors

In this section, we combine the graph sparsification and the partitioning algorithms using PageRank vectors to derive an improved partitioning algorithm.

An application of our sparsification algorithm by PageRank is the *balanced cut* problem. For a given graph  $G$ , we first use our sparsification algorithm to preprocess the graph. Then we apply the local partitioning algorithm using PageRank vectors [3, 4] on the sparsifier. Since the local partitioning algorithm is a subroutine for the balance cut problem, we obtain a balanced cut algorithm with an improved running time. Spielman and Teng [46] gave a local partitioning algorithm which, for a fixed value of  $\phi$ , gives a cut with approximation ratio  $O(\phi^{1/2} \log^{3/2} n)$  and of volume  $v_\phi$  in  $O(m \log^c(n)/\phi^{-5/3})$  time where  $v_\phi$  is the largest volume of the set with cheeger ratio  $\phi$ . Note that the constant  $c$  above is quite large [46]. In [3], PageRank vectors were used to derive a local partitioning algorithm with an improved running time  $(m + n\phi^{-1})O(\text{polylog}(n))$ . In [4], the running time was further reduced to  $O(m + n\phi^{-1/2})O(\text{polylog}(n))$  by preprocessing using sparsification algorithms in [9].

Given an undirected, weighted graph  $G = (V, E, w)$  with  $n$  vertices and  $m$  edges, we can apply algorithm **SparsifyApprGreen** as a preprocess procedure on graph  $G$  to get a sparsifier  $\tilde{G}$  with only  $O(n \log n/\epsilon^2)$  edges in time  $O(m \log^2 n/\phi)$  such that  $|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon h_G(S)$  for all  $S$  with  $h_G(S) \geq \phi$ . Then, we use the algorithm **PageRank-Partition** [3] on graph  $\tilde{G}$  instead of  $G$  for balanced cut problem. The algorithm

**PageRank-Partition** has two inputs including a parameter  $\phi$  and a graph with  $m$  edges. As stated in [3], the **PageRank-Partition** algorithm has expected running time  $O(m \log^4 m / \phi^2)$ . Furthermore, with high probability the **PageRank-Partition** algorithms was shown to be able to find a set  $S$ , if exists, such that  $\text{vol}(S) \geq \text{vol}(C)/2$  and  $h_G(S) \leq \phi$ . This can be summarized as follows:

**Theorem 9.** *Given an undirected, weighted graph  $G = (V, E, w)$  with  $n$  vertices and  $m$  edges, constant  $\phi > 0$ , and  $\epsilon > 0$ . With probability  $1/2$ , we can preprocess graph  $G$  in  $O(\frac{m \log^2 n}{\phi})$  time to obtain a sparse graph  $\tilde{G}$  with  $O(\frac{n \log n}{\epsilon^2})$  edges such that for all  $S \in V$  satisfying  $h_G(S) \geq \phi$ ,*

$$|h_{\tilde{G}}(S) - h_G(S)| \leq \epsilon h_G(S).$$

*Algorithm **PageRank-Partition** takes as inputs a parameter  $\phi$  and a graph  $\tilde{G}$  and has expected running time  $O(n \log^6 m / (\phi^2 \epsilon^2))$ . If there exists a set  $C$  with  $h_G(C) = O(\phi^2 / \log^2 n)$ , then with high probability the **PageRank-Partition** algorithm finds a set  $S$  such that  $\text{vol}(S) \geq \text{vol}(C)/2$  and  $h_G(S) \leq \phi$ .*

## References

1. D. Achlioptas, Database-friendly random projections, *PODS'01*, pages 274–281.
2. N. Alon, A. Andoni, T. Kaufman, K. Matulef, R. Rubinfeld, N. Xie, Testing  $k$ -wise and almost  $k$ -wise independence, *STOC 07*, 496–505
3. R. Andersen, F. Chung, and K. Lang, Local graph partitioning using pagerank vectors, *FOCS 06*, 475–486.
4. R. Andersen and Y. Peres, Finding sparse cuts locally using evolving sets, *STOC 09*, 235–244.
5. D. Achlioptas and F. McSherry, Fast computation of low rank matrix approximations, *STOC 01*, 611–618.
6. S. Arora, E. Hazan, and S. Kale,  $\Theta(\sqrt{\log n})$  approximation to sparsest cut in  $\tilde{O}(n^2)$  time, *FOCS 04*, 238–247.
7. S. Arora and S. Kale. A combinatorial, primal-dual approach to semidefinite programs, *STOC 07*, 227–236.
8. J. Batson, D. A. Spielman, and N. Srivastava, Twice-Ramanujan Sparsifiers, *STOC 09*, 255–262.
9. A. A. Bencúr and D. R. Karger. Approximating s-t minimum cuts in  $\tilde{O}(n^2)$  time, *STOC 96*, 47–55.
10. P. Berkhin, Bookmark-coloring approach to personalized pagerank computing, *Internet Mathematics*, **3(1)**, (2007), 41–62
11. S. Brin and L. Page, The anatomy of a large-scale hypertextual Web search engine, *Computer Networks and ISDN Systems*, **30(1-7)**, (1998), 107–117.
12. P. Drineas and R. Kannan, Pass efficient algorithms for approximating large matrices, *SODA 94*, 223–232.
13. F. Chung, *Spectral Graph Theory*, AMS Publication, 1997.
14. F. Chung and S.-T. Yau, Discrete Green's Functions, *Journal of Combinatorial Theory, Series A*, **91(1-2)**(2000), 191–214.
15. U. Feige and R. Krauthgamer, Finding and certifying a large hidden clique in a semi-random graph, *Random Structures and Algorithms*, **13**, (1998), 457–466.
16. A. Firat, S. Chatterjee, and M. Yilmaz, Genetic clustering of social networks using random walks, *Computational Statistics and Data Analysis*, **51(12)**, (2007), 6285–6294.
17. F. Fouss, A. Pirotte, J.-M. Renders, and M. Saerens, Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation, *IEEE Transactions on Knowledge and Data Engineering*, **19(3)** (2007), 355–369.
18. A. Frieze and R. Kannan, A new approach to the planted clique problem, *Foundations of Software Technology and Theoretical Computer Science (Bangalore) 2008*, (R. Hariharan, M. Mukund, V. Vinay eds.) 187–198.
19. A. Frieze, R. Kannan, and S. Vempala, Fast Monte-Carlo algorithms for finding low-rank approximations, *FOCS 98*, 370–378.
20. G. Green, An Essay on the Application of Mathematical Analysis to the Theories of Electricity and Magnetism, Nottingham, 1828.
21. J. Hastad, Clique is hard to approximate within  $n^{1-\epsilon}$ , *FOCS 97*, 627–636.
22. H. Haveliwala, Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search, *IEEE Trans. Knowl. Data Eng.*, **15(4)**, (2003), 784–796.
23. E. Hazan and R. Krauthgamer, How hard is it to approximate the best Nash equilibrium? *SODA 2009*, 720–727.
24. J. Hopcroft and D. Sheldon, Manipulation-resistant reputations using hitting time, *WAW 2007, LNCS 4853*, 68–81.
25. G. Jeh and J. Widom, SimRank: A measure of structural-context similarity, *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2002, 538–543.

26. G. Jeh and J. Widom, Scaling personalized web search, *Proceedings of the 12th World Wide Web Conference (WWW)* (2003), 271–279.
27. J. Kahn, J. H. Kim, L. Lovász, and V. H. Vu, The cover time, the blanket time, and the Matthews bound, *FOCS 00*, 467–475.
28. M. Jerrum, Large cliques elude the Metropolis process, *Random Structures and Algorithms*, 3 (1992), 347–359.
29. A. Juels and M. Peinado, Hiding cliques for cryptographic security, *SODA 98*, 677–684.
30. David R. Karger. Random sampling in cut, flow, and network design problems, *STOC 94*, 648–657.
31. David R. Karger. Using randomized sparsification to approximate minimum cuts, *SODA 94*, 424–432.
32. David R. Karger. Minimum cuts in near-linear time, *JACM*, 47 (2000), 46–76..
33. R. Karp, Reducibility among combinatorial problems, in *The complexity of computer computations*, (R. Miller and J. Thatcher eds.) Plenum Press, New York, (1972), 85–103.
34. R. Karp, The probabilistic analysis of some combinatorial search algorithms, in *Algorithms and Complexity: New Directions and Recent Results*, (J. F. Traub, ed.), Academic Press (1976), 1–19.
35. L. Kučera, Expected complexity of graph partitioning problems, *Discrete Applied Mathematics* 57, (1995), 193–212.
36. D. Liben-Nowell and J. Kleinberg, The link prediction problem for social networks, *Proceedings of the 12th International Conference on Information and Knowledge Management*, (2003), 556–559.
37. L. Lovász, Random walks on graphs: A survey, *Combinatorics, Paul Erdős is Eighty* 2 (1993), 1–46.
38. L. Lovász and M. Simonovits, The mixing rate of Markov chains, an isoperimetric inequality, and computing the volume, *FOCS 90*, 346–354.
39. M. Mihail, Conductance and convergence of markov chains—a combinatorial treatment of expanders, *FOCS 89*, 526–531.
40. L. Minder and D. Vilenchik, Small clique detection and approximate Nash equilibria, *Approx and Random 2009, LNCS 5687* 673–685.
41. L. Orecchia, L. J. Schulman, U. V. Vazirani, and N. K. Vishnoi. On partitioning graphs via single commodity flows, *STOC 08*, 461–470.
42. L. Page, S. Brin, R. Motwani, and T. Winograd, The pagerank citation ranking: Bringing order to the web, Technical report, Stanford Digital Library Technologies Project, 1998.
43. P. A. Pevzner and S. Sze, Combinatorial approaches to finding subtle signals in DNA sequences, *Proceedings of the Eight International Conference on Intelligent Systems for Molecular Biology*, (2000), 269–278.
44. M. Rudelson and R. Vershynin, Sampling from large matrices: An approach through geometric functional analysis, *Journal of the ACM*, 54(4) (2007), Article 21.
45. D. A. Spielman and S.-H. Teng. Spectral partitioning works: Planar graphs and finite element meshes, *FOCS 96*, 96–105.
46. D. A. Spielman and S.-H. Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. *STOC 04*, 81–90.
47. D. A. Spielman and S.-H. Teng. Nearly-linear time algorithms for preconditioning and solving symmetric, diagonally dominant linear systems. Available at <http://www.arxiv.org/abs/cs.NA/0607105>, 2006.
48. D. A. Spielman and N. Srivastava, Graph sparsification by effective resistances. *STOC 2008*, 563–568.