

Locating highly connected clusters in large networks with HyperLogLog counters

Lotte Weedage*, Nelly Litvak*, Clara Stegehuis*

January 13, 2021

Abstract

In this paper we introduce a new method to locate highly connected clusters in a network. Our proposed approach adapts the HYPERBALL algorithm [8] to localize regions with a high density of small subgraph patterns in large graphs in a memory-efficient manner. We use this method to evaluate three measures of subgraph connectivity: conductance, the number of triangles, and transitivity. We demonstrate that our algorithm, applied to these measures, helps to identify clustered regions in graphs, and provides good seed sets for community detection algorithms such as PageRank-Nibble. We analytically obtain the performance guarantees of our new algorithms, and demonstrate their effectiveness in a series of numerical experiments on synthetic and real-world networks. hyperloglog, probabilistic counting, network clustering

1 Introduction

Networks describe the connections between pairs of objects. Examples of networks are ubiquitous and include social networks, the Internet or communication networks. While these examples are very different from an application point of view, they share many characteristics. For example, in many real-world networks, objects have the tendency to cluster together in groups. The task of finding these densely connected spots, or clusters, in the network, has been a subject of vast research.

A quantity that is commonly used to measure the quality of clusters in a graph representation of a network, is the conductance of a cluster. The conductance is based on the *min-cut* [20], and gives a ratio of the number of edges connected to nodes outside the cluster relative to the number of edges inside the cluster. The use of conductance as a measure to find communities is one of the most useful and important cut-based methods that exists [18].

Another measure that can find clustered groups of nodes in a graph is the number of triangles, because a triangle in a graph is the most clustered subgraph consisting of three nodes. For any subgraph in a network, one can measure its quality as a cluster by counting the number of triangles in this subgraph. Another option is to find its *transitivity*, which is the ratio of the number of triangles versus the number of wedges in the subgraph, and therefore it tells us how clustered the graph is. Moreover, the task of finding triangles itself has interesting applications in for example, biological networks [21], spam detection [4] or link recommendations [22].

*Department of Electrical Engineering, Mathematics and Computer Science, University of Twente, the Netherlands

Finding dense parts of a graph, that have a low conductance, a high number of triangles, or a high transitivity, is a computationally demanding task because real-world networks often contain millions or even billions of nodes. In this paper we propose new methods to efficiently compute three measures of clustering – the conductance, the number of triangles, and the transitivity – for the ball subgraphs $\mathcal{B}_r(v)$: the induced subgraphs that contain all nodes within graph distance r of node v . The identified ball subgraphs of low conductance or high transitivity reveal locations of dense areas in the network and can be used in community detection, for example, as *seed sets* of other more time- and memory-consuming algorithms such as PAGERANK NIBBLE [1] or the MULTI WALKER CHAIN model [6].

Our proposed algorithms for computing conductance and transitivity use probabilistic HyperLogLog counters to estimate the number of edges, wedges and triangles in ball subgraphs. This class of randomized algorithms stems from the HYPERLOGLOG algorithm [10] for counting the number of distinct elements in large streams of data, such as the number of unique visitors on a web page or the number of different genomes in biological data. These counters give an accurate estimate of large cardinalities, and moreover are memory-efficient. Boldi and Vigna [7] have successfully adapted the HYPERLOGLOG algorithm to the networks context by developing the HYPERBALL algorithm, which counts the number of nodes in a ball subgraph $\mathcal{B}_r(v)$ for every node v and every radius r . They used this to approximate centrality measures in large graphs and to find the distribution of distances between pairs of nodes in a network of Facebook users [2]. This idea of counting nodes in ball subgraphs has also been extended to counting distinct edges in ball subgraphs [13] or in a stream of edges [25]. In this paper, we demonstrate that the potential of HyperLogLog-type counters on graphs is greater than only counting nodes and edges, but can extend to counting other patterns in networks, and can be used to approximate popular measures of clustering.

The main contribution of this paper is in designing memory-efficient HyperLogLog-based algorithms for computing the conductance, the number of triangles, and the transitivity in ball subgraphs. We analytically derive accurate error bounds for these algorithms, and empirically confirm their high performance. Moreover, we demonstrate applications of our methods to community detection in synthetic and real-world networks. Our results show that the identified highly clustered ball subgraphs perform very well as seed sets of the PR-NIBBLE algorithm [1], improving on previously used benchmarks.

The structure of this paper is as follows. In Section 2, we provide a brief recap on algorithms for counting nodes and edges in graphs using HyperLogLog counters. In Section 3, we extend these algorithms to other patterns in graphs, and we present our new methods for approximating the conductance, the number of triangles, and the transitivity in ball subgraphs. In Section 4, we analytically derive accurate error bounds of the estimators for the conductance, triangle count and transitivity. In Section 5, we experimentally evaluate the performance of our algorithms on a number of synthetic networks. In Section 6, we demonstrate application of our methods to community detection. We conclude in Section 7 with discussion.

2 Preliminaries: HYPERLOGLOG, HYPERBALL, and HYPEREDGE-BALL

2.1 HYPERLOGLOG algorithm

The HYPERLOGLOG algorithm is a probabilistic counting technique that estimates the number of distinct elements of a large dataset, called the dataset cardinality. While a naive deterministic approach requires storage of all distinct elements observed so far, the randomized HYPERLOGLOG algorithm is extremely memory-efficient, yet delivers accurate cardinality estimations.

We will now briefly outline the idea of the HYPERLOGLOG algorithm [10]. The input of the algorithm is

a multiset \mathcal{M} , a stream of data items that are read in order of occurrence. The algorithm uses a hash function $h : \mathcal{M} \rightarrow \{0, 1\}^\infty$ that assigns a binary string to every element of \mathcal{M} . The hash function is deterministic in the sense that it assigns exactly the same value to identical elements of \mathcal{M} . However, h is constructed in such a way that its bits can be assumed to be independent Bernoulli random variables with probability $1/2$ of 0 and 1. Then one can use the principle of *bit-pattern observables*: for example, in order to encounter the pattern 00001 at the beginning of a string, one needs to observe, on average, 32 different items. For such estimate, the algorithm needs to store in memory only ‘how rare’ the ‘most rare’ observed binary sequence is, for example, the maximal number of zeros observed so far at the beginning of a binary sequence. The name HYPERLOGLOG refers to this extremely low, double-logarithmic, memory requirements. To obtain accurate estimates, the algorithm uses registers. That is, the first b bits of h are used for identifying one of the $p = 2^b$ registers, and the string after that is used to compute the cardinality estimate in this register. The algorithm initialises an empty counter with $p = 2^b$ registers, where every register corresponds to an entry of the counter. The more registers we use, the more precise the cardinality estimate will be. More precisely, the HYPERLOGLOG algorithm returns the following estimate E of the number of unique elements of multiset \mathcal{M} :

$$E := \frac{\alpha_p p^2}{\sum_{j=1}^p 2^{-M[j]}}, \quad \text{with } \alpha_p := \left(p \int_0^\infty \left(\log_2 \left(\frac{2+u}{1+u} \right) \right)^p du \right)^{-1}, \quad (1)$$

where $M[j]$ is the cardinality estimate of register j . The pseudocode of the HYPERLOGLOG algorithm is provided in Algorithm A.1 in Appendix A. Overall, the HYPERLOGLOG algorithm uses $(1+o(1))p \log \log(n/p)$ bits of space [10] for a set of cardinality n , making it an extremely memory-efficient algorithm to estimate cardinalities of large sets.

The expectation and the variance of E are given in Theorem 1 of [10], and will be used later in the paper for obtaining performance guarantees of our algorithms:

Definition 1 (Ideal multiset [10]). *An ideal multiset of cardinality n is a sequence obtained by arbitrary replications and permutations applied to n uniform identically distributed random variables over the real interval $[0, 1]$.*

Theorem 1 (from [10]). *Let the algorithm HYPERLOGLOG be applied to an ideal multiset of (unknown) cardinality n , using $p \geq 3$ registers, and let E be the resulting cardinality estimate.*

(i) *The estimate E is asymptotically almost unbiased in the sense that, as $n \rightarrow \infty$,*

$$\frac{1}{n} \mathbb{E}(E) = 1 + \delta_1(n) + o(1), \quad \text{where } |\delta_1(n)| < 5 \cdot 10^{-5} \text{ as soon as } p \geq 2^4. \quad (2)$$

(ii) *The standard error defined as $\frac{1}{n} \sqrt{\text{Var}(E)}$ satisfies as $n \rightarrow \infty$,*

$$\frac{1}{n} \sqrt{\text{Var}(E)} = \frac{\beta_p}{\sqrt{p}} + \delta_2(n) + o(1), \quad \text{where } |\delta_2(n)| < 5 \cdot 10^{-4} \text{ as soon as } p \geq 2^4, \quad (3)$$

the constants β_p being bounded, with $\beta_{16} = 1.106, \beta_{32} = 1.070, \beta_{64} = 1.054, \beta_{128} = 1.046$, and $\beta_\infty = \sqrt{3 \log(2) - 1} \approx 1.03896$.

2.2 HYPERBALL algorithm

The HYPERBALL algorithm, introduced in [8], estimates the size of the ball consisting of nodes within graph distance r around a center node using HyperLogLog counters. The algorithm is an adaptation of the HYPERANF algorithm [7], which is based on the fact that the nodeball around node v with radius r , $\mathcal{B}_r(v)$, can be found iteratively:

Definition 2 (Nodeball). *The nodeball $\mathcal{B}_r(v)$ consists of every node in a ball of radius r around node v . Define $\mathcal{B}_0(v) = \{v\}$. For $r > 1$, define:*

$$\mathcal{B}_r(v) = \bigcup_{w:\{v,w\} \in E} \mathcal{B}_{r-1}(v) \cup \mathcal{B}_{r-1}(w). \quad (4)$$

The HYPERBALL algorithm uses one HyperLogLog counter per node and for each iteration, the counters of the neighbours of this node are added to the node's own counter. After each iteration r , the size of this counter is calculated, which equals the estimator of $|\mathcal{B}_{r+1}(v)|$. By Theorem 1, every estimator is almost unbiased and has a relative error of at most β_p/\sqrt{p} , for $\beta_p < 1.046$ as soon as every HyperLogLog counter has $p > 128$ registers. This algorithm particularly excels at its small memory usage and the fact that the size of nodeballs around all nodes are found simultaneously. The pseudocode of the HYPERBALL algorithm is provided in Algorithm A.2 in Appendix A.

2.3 HYPEREDGEBALL algorithm

The HYPERBALL algorithm naturally extends from counting nodes to counting edges that can be reached within radius r around a node [13]. The corresponding edgeballs are defined as follows:

Definition 3 (Edgeball). *The edgeball $\mathcal{E}_0(v)$ consists of every edge incident to node v . Then, for $r > 1$:*

$$\mathcal{E}_r(v) = \bigcup_{w:\{v,w\} \in E} \mathcal{E}_{r-1}(v) \cup \mathcal{E}_{r-1}(w). \quad (5)$$

Equivalently, we can rewrite (5) as

$$\mathcal{E}_r(v) = \{\{x, y\} : x \in \mathcal{B}_r(v)\}. \quad (6)$$

The HYPEREDGEBALL algorithm gives an approximation of the number of edges around a node after r iterations, $|\mathcal{E}_{r+1}(v)|$. Compared to the HYPERBALL algorithm, the only change is in the initialisation phase (Algorithm A.2), since we now count edges instead of nodes. This new initialisation is formalised in Algorithm 1.

3 Algorithms

We will now show that we can easily extend HyperBall-type algorithms to counting arbitrary, more complex patterns than edges or nodes. In particular, we will show that we can apply HyperBall-type algorithms to approximate three important clustering measures: conductance, triangles, and transitivity.

Algorithm 1 The HYPEREDGEBALL algorithm. The ADD and SIZE functions of Algorithm A.1 and the COUNTBALL function of Algorithm A.2 are used to find the edgeball estimators.

```

1:  $c[-]$ , an array of  $n$  empty HyperLogLog counters
2:
3: for each  $v \in V$  do
4:   for each  $e \in \{\{v, w\} \in E\}$  do
5:     ADD( $c[v], w$ )
6:   end for
7:   write  $\langle v, c[v] \rangle$  to disk
8:   return SIZE( $c[v]$ ), which estimates  $|\mathcal{E}_0(v)|$ 
9: end for
10:
11:  $(\widehat{\mathcal{E}}_r)_{r \geq 1} = \text{COUNTBALL}(c)$ 

```

3.1 Conductance

The first quantity of interest that we investigate is conductance. The conductance of a graph G is defined as follows:

Definition 4 (Conductance). *For a graph $G = (V, E)$ with $n = |V|$ and $m = |E|$, the conductance of a subgraph $S \subset G$ is:*

$$\phi(S) = \frac{|\delta(S)|}{\min(\text{vol}(S), 2m - \text{vol}(S))}, \quad (7)$$

where $\delta(S) = \{\{x, y\} \in E \mid x \in S, y \notin S\}$ is the boundary of S , and $\text{vol}(S)$ is the volume of S equal to the sum of the degrees of the nodes in subgraph S .

In the rest of this paper we use a simplified definition of the conductance, as we assume that the graphs that we analyse are non-empty and the volume of the subgraphs that we analyze is smaller than the volume of its complement:

$$\phi(S) = \frac{|\delta(S)|}{\text{vol}(S)}. \quad (8)$$

Next, we wish to provide a memory-efficient way of estimating conductance of ball subgraphs. To this end, it seems natural to combine HYPERBALL and HYPEREDGEBALL from Section 2. However, this turns out to be insufficient because the low memory usage implies that we cannot identify the edges of $\delta(\mathcal{B}_r(v))$. To overcome this problem, we transform our (undirected) graph into a directed variant of the graph where every undirected edge becomes two directed edges, and we introduce directed edgeballs:

Definition 5 (Out-edgeball). *The out-edgeball with radius r around node v is defined as follows:*

$$\mathcal{E}_r^-(v) = \{(x, y) \in E \mid x \in \mathcal{B}_r(v)\}. \quad (9)$$

Similarly, we can define an in-edgeball:

Definition 6 (In-edgeball). *The in-edgeball with radius r around node v is defined as follows:*

$$\mathcal{E}_r^+(v) = \{(x, y) \in E \mid y \in \mathcal{B}_r(v)\}. \quad (10)$$

The different kinds of edgeballs are illustrated in Figure 1. For our purposes, it is important to notice that when an edge is on the boundary between a node inside and a node outside the ball of radius r , this edge belongs only to the out-edgeball, while all other edges, that have both endpoints of the edge are in the nodeball $\mathcal{B}_r(v)$, belong to both in- and out-edgeball. This is exactly why the directed edgeballs are helpful for estimating conductance. This is formally stated in the following theorem:

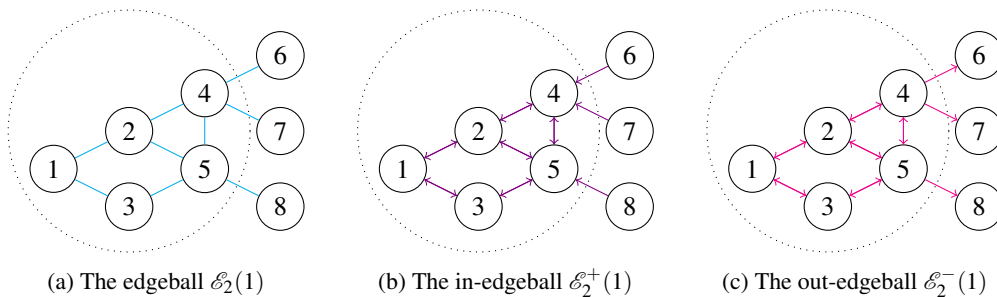


Figure 1: Three methods for counting edges in an undirected graph.

Theorem 2. For an undirected ball subgraph $S_r(v)$ it holds that:

$$|\delta(\mathcal{B}_r(v))| = 2|\mathcal{E}_r(v)| - |\mathcal{E}_r^-(v)|, \quad (11)$$

$$\text{vol}(\mathcal{B}_r(v)) = |\mathcal{E}_r^-(v)|. \quad (12)$$

Proof. Denote by $\mathbf{1}\{A\}$ the indicator of A , and let E' be the set of directed edges obtained from the edge set E by replacing each undirected edge $\{x,y\}$ by two directed edges (x,y) and (y,x) . By definition of the volume of a set of nodes, we obtain

$$\text{vol}(\mathcal{B}_r(v)) = \sum_{x \in \mathcal{B}_r(v)} \sum_{y \in V} \mathbf{1}\{\{x,y\} \in E\} = \sum_{x \in \mathcal{B}_r(v)} \sum_{y \in V} \mathbf{1}\{(x,y) \in E'\} = |\mathcal{E}_r^-(v)|,$$

which proves (12). Next, using (12), we write

$$\begin{aligned} |\delta(\mathcal{B}_r(v))| &= \sum_{x \in \mathcal{B}_r(v)} \sum_{y \notin \mathcal{B}_r(v)} \mathbf{1}\{\{x,y\} \in E\} \\ &= \sum_{x \in \mathcal{B}_r(v)} \sum_{y \notin \mathcal{B}_r(v)} \mathbf{1}\{\{x,y\} \in E\} + \text{vol}(\mathcal{B}_r(v)) - \text{vol}(\mathcal{B}_r(v)) \\ &= \sum_{x \in \mathcal{B}_r(v)} \sum_{y \notin \mathcal{B}_r(v)} \mathbf{1}\{\{x,y\} \in E\} + \sum_{x \in \mathcal{B}_r(v)} \sum_{y \in V} \mathbf{1}\{\{x,y\} \in E\} - |\mathcal{E}_r^-(v)| \\ &= 2 \sum_{x \in \mathcal{B}_r(v)} \sum_{y \notin \mathcal{B}_r(v)} \mathbf{1}\{\{x,y\} \in E\} + \sum_{x \in \mathcal{B}_r(v)} \sum_{y \in \mathcal{B}_r(v)} \mathbf{1}\{\{x,y\} \in E\} - |\mathcal{E}_r^-(v)|. \end{aligned}$$

The last expression equals to the right-hand side of (11) by the definition of $\mathcal{E}_r(v)$ and the fact that in the second term each undirected edge is counted twice. \square

Theorem 2 suggests the following estimator $\hat{\phi}(\mathcal{B}_r(v))$ of the conductance of $\mathcal{B}_r(v)$:

$$\hat{\phi}(\mathcal{B}_r(v)) = \frac{|\delta(\mathcal{B}_r(v))|}{\text{vol}(\mathcal{B}_r(v))} = \frac{2\widehat{|\mathcal{E}_r(v)|} - \widehat{|\mathcal{E}_r^-(v)|}}{\widehat{|\mathcal{E}_r^-(v)|}} = 2 \frac{\widehat{|\mathcal{E}_r(v)|}}{\widehat{|\mathcal{E}_r^-(v)|}} - 1, \quad (13)$$

where $|\widehat{\mathcal{E}_r(v)}|$ and $|\widehat{\mathcal{E}_r^-(v)}|$ are the estimates of $|\mathcal{E}_r(v)|$ and $|\mathcal{E}_r^-(v)|$, respectively.

Algorithm 2 The directed HYPEREDGEBALL algorithm. The ADD and SIZE functions of Algorithm A.1 and the COUNTBALL function of Algorithm A.2 is used to find the out-edgeball estimators.

```

1:  $c[-]$ , an array of  $n$  empty HyperLogLog counters
2:
3: for each  $v \in V$  do
4:   for each  $e \in \{(v, w) \in E'\}$  do
5:     ADD( $c[v], w$ )
6:   end for
7:   write  $\langle v, c[v] \rangle$  to disk
8:   return SIZE( $c[v]$ ), which estimates  $|\mathcal{E}_0^-(v)|$ 
9: end for
10:
11:  $(|\widehat{\mathcal{E}_r^-}|)_{r \geq 1} = \text{COUNTBALL}(c)$ 

```

3.2 Triangles and wedges

We now present our algorithm that counts the number of triangles within a ball of radius r around node v . We denote this number by $\Delta_r(v)$, and its estimator by $\hat{\Delta}_r(v)$. The idea is to obtain $\hat{\Delta}_r(v)$ using the same HYPERBALL algorithm as for counting edges or nodes, but we initialise the counter with triangles instead of nodes or edges. For this initialisation, we assign a unique hash value to each triangle in the graph. This can be done using algorithms for exact triangle counting such as *compact-forward* [15] or *edge-iterator* [3]. In Algorithm 3, we give an example of how this can be implemented.

Now denote by $w_r(v)$ the number of wedges in $\mathcal{B}_r(v)$. Since wedges are open triangles, an estimator $\hat{w}_r(v)$ of $w_r(v)$ can be found in exactly the same way as $\hat{\Delta}_r(v)$, but in the initialisation we add a wedge $\{i, v, j\}$ to a counter if $i, j \in V$ are neighbours of v . (Note that in line 7 of Algorithm 3 we verify that i and j are neighbours; this is needed for the initialisation of the counter of triangles).

3.3 Extension to counting of arbitrary induced subgraphs

The algorithms above for counting triangles and wedges easily extend to counting arbitrary induced connected subgraphs, called *graphlets*, within balls $\mathcal{B}_r(v)$. For that, in the initialisation phase, we need to count the graphlets that involve node v for all $v \in V$, and assign a unique hash value to each of these graphlets. After that, we can run the HYPERBALL algorithm to count graphlets in ball subgraphs. This approach potentially can indicate parts of networks with unusual quantities of particular graphlets. However, applying HYPERBALL to general subgraph patterns also comes with important difficulties. First, HYPERBALL counts nodes/edges/graphlets in ball subgraphs, and therefore this approach cannot be easily extended to counting graphlets in subgraphs of any other form. Second, the initialisation phase presents a computational bottleneck because assigning a unique hash value to each graphlet supersedes the computationally demanding task of exact graphlet counting.

Algorithm 3 Triangle and wedge ball algorithm. The ADD and SIZE functions of Algorithm A.1 and the COUNTBALL function of Algorithm A.2 is used to find the triangle or wedge cardinality estimates.

```

1:  $c[-]$ , an array of  $n$  HyperLogLog counters
2: triangle, wedge: booleans that express whether triangles or wedges are counted
3:
4: for each  $v \in V$  do
5:   for each  $(i, v) \in E$  do
6:     for each  $(j, v) \in E$  do
7:       if  $(i, j) \in E$  and triangle then
8:         ADD( $c[v], (v, i, j)$ )
9:       else if wedge then
10:        ADD( $c[v], (v, i, j)$ )
11:      end if
12:    end for
13:  end for
14:  write  $\langle v, c[v] \rangle$  to disk
15:  return SIZE( $c[v]$ ), which estimates  $|\Delta_0(v)|$  or  $|w_0(v)|$ 
16: end for
17:
18: if triangle then
19:    $(\hat{\Delta}_r)_{r \geq 1} = \text{COUNTBALL}(c)$ 
20: else if wedge then
21:    $(\hat{t}_r)_{r \geq 1} = \text{COUNTBALL}(c)$ 
22: end if

```

4 Error bounds

4.1 Error bounds for the estimator of conductance

We now introduce Theorem 3 and 5 that give a lower and upper bound for the conductance estimator $\hat{\phi}(\mathcal{B}_r(v))$ from (13) based on Chebyshev's inequality and Vysochanskij-Petunin's inequality [24].

Theorem 3 (Chebyshev bound for the conductance estimator). *For all $v \in V$, $r \geq 1$, the conductance estimator $\hat{\phi}(\mathcal{B}_r(v))$, as defined in (13), satisfies:*

$$P \left[\hat{\phi}(\mathcal{B}_r(v)) \in \left(\frac{1 - \varepsilon}{1 + \gamma} \cdot \phi(S_r(v)), \frac{1 + \varepsilon}{1 - \gamma} \cdot \phi(S_r(v)) \right) \right] \geq 1 - \eta^2 \left(\frac{|\mathcal{E}_r^+(v)|^2}{p_1^2} + \frac{|\mathcal{E}_r^-(v)|^2}{p_2^2} \right), \quad (14)$$

with $\varepsilon = \frac{p_1}{|\mathcal{E}_r^+(v)|} + \delta_1 + o_1(|\mathcal{E}_r^+(v)|)$, $\gamma = \frac{p_2}{|\mathcal{E}_r^-(v)|} + \delta_1 + o_2(|\mathcal{E}_r^-(v)|)$ and $p_1, p_2 > 0$, where $\delta_1 = 5 \cdot 10^{-5}$, and $o_1, o_2 = o(1)$, as their argument goes to infinity.

The proof of this theorem is based on Theorem 1, and is given in Appendix B.

When our estimators $|\widehat{\mathcal{E}_r^+(v)}|$ and $|\widehat{\mathcal{E}_r^-(v)}|$ follow a unimodal distribution, we can also use Vysochanskij-Petunin's (VP) inequality [24]:

Algorithm 4 Graphlet ball algorithm. The ADD and SIZE functions of Algorithm A.1 and the COUNTBALL function of Algorithm A.2 is used to find the graphlet cardinality estimates.

```

1:  $c[-]$ , an array of  $n$  HyperLogLog counters
2:
3: for each  $v \in V$  do
4:   if  $v \in \text{graphlet}$  then
5:     ADD( $c[v]$ , graphlet)
6:     write  $\langle v, c[v] \rangle$  to disk
7:     return SIZE( $c[v]$ ), which estimates the number of graphlets in  $\mathcal{B}_1(v)$ .
8:   end if
9: end for
10:
11:  $(\text{graphlet estimate}_r)_{r \geq 1} = \text{COUNTBALL}(c)$ 

```

Theorem 4 (Vysochanskij-Petunin’s inequality). *Assume that the random variable X has a unimodal distribution with finite mean $E(X)$ and variance $\text{Var}(X) = \sigma^2$. Then, for $\lambda/\sigma > \sqrt{8/3}$,*

$$P(|X - E(X)| \geq \lambda) \leq \frac{4\sigma^2}{9\lambda^2}. \quad (15)$$

By using this inequality instead of Chebyshev’s inequality, we can obtain tighter error bounds than the ones in Theorem 3, given in the next theorem.

Theorem 5 (Vysochanskij-Petunin bound for the conductance estimator). *If $|\widehat{\mathcal{E}}_r(v)|$ and $|\widehat{\mathcal{E}}_r^-(v)|$ have a unimodal distribution, then for all $v \in V$, $r \geq 1$, the conductance estimator $\hat{\phi}(\mathcal{B}_r(v))$, as defined in (13), satisfies:*

$$P\left[\hat{\phi}(\mathcal{B}_r(v)) \in \left(\frac{1-\varepsilon}{1+\gamma} \cdot \phi(S_r(v)), \frac{1+\varepsilon}{1-\gamma} \cdot \phi(S_r(v))\right)\right] \geq 1 - \frac{4}{9}\eta^2 \left(\frac{|\mathcal{E}_r(v)|^2}{\lambda_1^2} + \frac{|\mathcal{E}_r^-(v)|^2}{\lambda_2^2}\right), \quad (16)$$

with $\lambda_1 > \sqrt{8/3 \cdot \text{Var}(|\mathcal{E}_r(v)|)}$, $\lambda_2 > \sqrt{8/3 \cdot \text{Var}(|\mathcal{E}_r^-(v)|)}$, $\varepsilon = \frac{\lambda_1}{|\mathcal{E}_r(v)|} + \delta_1 + o_1(|\mathcal{E}_r(v)|)$ and $\gamma = \frac{\lambda_2}{|\mathcal{E}_r^-(v)|} + \delta_1 + o_2(|\mathcal{E}_r^-(v)|)$, where δ_1, o_1 and o_2 are as in Theorem 3.

The proof of this theorem is identical to the proof of Theorem 3, but it uses the Vysochanskij-Petunin inequality instead of the Chebyshev’s inequality. In Section 5 we will show that the VP bound indeed holds in our numerical experiments by using a statistical test of unimodality [12]. In future research it will be interesting to identify general conditions under which the HyperLogLog counters produce estimators with a unimodal distribution.

4.2 Error bounds of the estimator for the triangle count

We again use the Chebyshev inequality in order to find a lower and upper error bound for the triangle count estimator $\hat{\Delta}_r(v)$ from Algorithm 3:

Theorem 6 (Chebyshev bound for the triangle count estimator). *For all $v \in V$, $r \geq 1$, the triangle estimator $\hat{\Delta}_r(v)$ satisfies:*

$$P\left[\hat{\Delta}_r(v) \in \left(\mathbb{E}(\hat{\Delta}_r(v)) - a, \mathbb{E}(\hat{\Delta}_r(v)) + a\right)\right] \geq 1 - \frac{\eta^2 \Delta_r(v)^2}{a^2}, \quad (17)$$

for $a > 0$ and $\eta = \frac{\beta_p}{\sqrt{p}} + \delta_2 + o_1(\Delta_r(v))$, where $\delta_2 = 5 \cdot 10^{-4}$ and $o_1 = o(1)$, as its argument goes to infinity.

The proof of this theorem is again based on Theorem 1, and is given in Appendix B.

The Vysochanskij-Petunin inequality [24], which holds whenever an estimator has a unimodal distribution, can also be used in order to get a slightly tighter error bound:

Theorem 7 (Vysochanskij-Petunin bound for the triangle estimator). *If $\hat{\Delta}_r(v)$ has a unimodal distribution, then for all $v \in V$, $r \geq 1$, the triangle estimator $\hat{\Delta}_r(v)$ satisfies:*

$$P\left(\hat{\Delta}_r(v) \in \left(\mathbb{E}(\hat{\Delta}_r(v)) - \lambda, \mathbb{E}(\hat{\Delta}_r(v)) + \lambda\right)\right) \geq 1 - \frac{4\eta^2 \Delta_r(v)^2}{9\lambda^2}, \quad (18)$$

for $\lambda > \sqrt{8/3 \cdot \text{Var}(\Delta_r(v))}$ and η as in Theorem 6.

The proof of Theorem 7 goes in the same way as the proof of Theorem 6, but with the Vysochanskij-Petunin inequality instead of Chebyshev's inequality.

4.3 Error bounds for transitivity

The transitivity of a graph G is defined as follows:

Definition 7 (Transitivity). *The transitivity of a graph G equals to*

$$t(G) = \frac{3\Delta(G)}{w(G)}, \quad (19)$$

where $w(G)$ is the number of wedges, and $\Delta(G)$ the number of triangles in G .

In order to find the transitivity of ball subgraphs, we need to find the number of wedges, $|w(\mathcal{B}_r(v))|$, in these ball subgraphs, which we obtain by using Algorithm 3.

4.3.1 Error bounds of the transitivity estimator

We can find the Chebyshev and Vysochanskij-Petunin error bounds of our estimator for transitivity similarly to Theorem 6 and 7. The transitivity estimator is

$$\hat{t}(\mathcal{B}_r(v)) := \frac{3\hat{\Delta}_r(v)}{\hat{w}_r(v)}. \quad (20)$$

The expectation and variance of our estimators $\hat{\Delta}_r(v)$ and $\hat{w}_r(v)$ can be obtained from Theorem 1, which results in the following Chebyshev error bound of the transitivity estimator:

Theorem 8 (Chebyshev bound for the transitivity estimator). *For $v \in V, r \geq 1$, the transitivity estimator $\hat{t}(\mathcal{B}_r(v))$, as defined in (20), satisfies:*

$$P \left[\hat{t}(\mathcal{B}_r(v)) \in \left(\frac{1-\varepsilon}{1+\gamma} \cdot t(\mathcal{B}_r(v)), \frac{1+\varepsilon}{1-\gamma} \cdot t(\mathcal{B}_r(v)) \right) \right] \geq 1 - \eta^2 \left(\frac{(\Delta_r(v))^2}{p_1^2} + \frac{(w_r(v))^2}{p_2^2} \right), \quad (21)$$

for $p_1, p_2 > 0$ and $\varepsilon = \frac{p_1}{\Delta_r(v)} + \delta_1 + o_1(\Delta_r(v))$, $\gamma = \frac{p_2}{w_r(v)} + \delta_1 + o_2(w_r(v))$, with $\delta_1 = 5 \cdot 10^{-5}$ and $o_1, o_2 = o(1)$ as their argument goes to infinity.

When this transitivity estimate has a unimodal distribution amongst the nodes, we can again use the Vysochanskij-Petunin inequality in order to get tighter error bounds:

Theorem 9 (Vysochanskij-Petunin bound for the transitivity estimator). *If $\hat{\Delta}_r(v)$ and $\hat{w}_r(v)$ have a unimodal distribution, then for all $v \in V, r \geq 1$, the transitivity estimator $\hat{t}(\mathcal{B}_r(v))$, as defined in (20), satisfies:*

$$P \left[\hat{t}(\mathcal{B}_r(v)) \in \left(\frac{1-\varepsilon}{1+\gamma} \cdot t(\mathcal{B}_r(v)), \frac{1+\varepsilon}{1-\gamma} \cdot t(\mathcal{B}_r(v)) \right) \right] \geq 1 - \frac{4}{9} \eta^2 \left(\frac{\Delta_r(v)^2}{\lambda_1^2} + \frac{w_r(v)^2}{\lambda_2^2} \right), \quad (22)$$

with $\lambda_1 > \sqrt{8/3 \cdot \text{Var}(\Delta_r(v))}$, $\lambda_2 > \sqrt{8/3 \cdot \text{Var}(w_r(v))}$, $\varepsilon = \frac{\lambda_1}{\Delta_r(v)} + \delta_1 + o_1(\Delta_r(v))$ and $\gamma = \frac{\lambda_2}{w_r(v)} + \delta_1 + o_2(w_r(v))$, with δ_1, o_1 and o_2 as in Theorem 8.

5 Performance

To evaluate the performance of our algorithms, we first run a series of experiments on artificial LFR graphs [14]. In LFR graphs, the nodes are divided into pre-defined communities, and one can choose a *mixing parameter* $\mu \in [0, 1]$, which is the probability that an edge emanating from a node connects to a node in outside of its community. The LFR model involves a number of other parameters: the minimum and maximum community size ($|C_{min}|, |C_{max}|$), the average and maximum degree (\bar{d}, d_{max}), the power-law exponent of the inverse cumulative distribution of the node degrees (τ_1), and the power-law exponent of the inverse cumulative distribution of community sizes (τ_2). We have used the LFR graph generator of NetworkX 2.4 in Python 3.6.9 to create LFR graphs with three different sets of parameters as shown in Table 1. We have used graphs with 1000 and 5000 nodes because in these small graphs we can find the exact number of edges, directed edges, wedges and triangles in all ball graphs, and compare the performance of our algorithms to these exact results.

Graph	$n = V $	τ_1	τ_2	$ C_{min} $	$ C_{max} $	\bar{d}	d_{max}	μ
LFR-1	1000	2	3	10	50	10	50	0.3
LFR-2	1000	2	3	20	100	20	100	0.3
LFR-3	5000	2	3	10	50	10	50	0.3

Table 1: Parameters for the generated LFR-graphs

For the experiments on LFR graphs, we have used a mixing parameter $\mu = 0.3$, since this was the smallest mixing parameter with no isolated nodes in every generated graph. We have investigated the conductance, number of triangles and transitivity in ball subgraphs of radii 1 and 2. A ball larger than this radius consists of a large part of the entire graph.

5.0.1 Conductance

Figure 2 shows the exact and the estimated conductance in a LFR-3 graph in ball subgraphs of radius 1 (Figure 2a) and radius 2 (Figure 2b). On the horizontal axis, the nodes v are arranged in the order of ascending conductance of $\mathcal{B}_r(v)$, $r = 1, 2$. The DIP test for unimodality [12] gives a small p -value of 0.0085, therefore we can reasonably assume that the conductance estimator is unimodal and apply the Vysochanskij-Petunin error bounds (Theorem 5). Figure 2 shows that the Vysochanskij-Petunin bounds are tight and represent well the 95%-margin of the estimation.

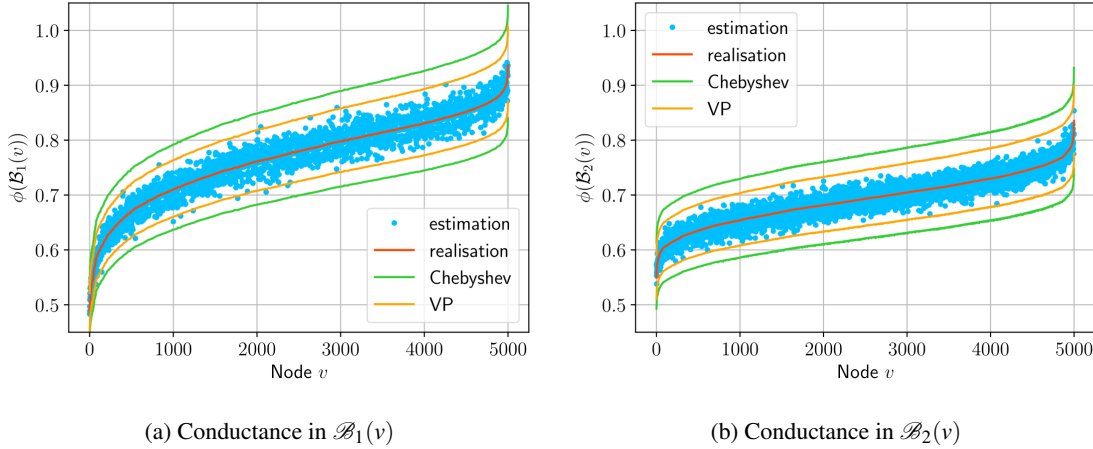


Figure 2: Estimated and exact conductance in the LFR-3 graph with $p = 2^{14}$ registers. The nodes v are sorted by realised conductance of $\mathcal{B}_r(v)$ in ascending order.

In order to investigate how the precision improves with increasing the number of registers, We have experimented with different numbers of registers in the LFR-1 graph. The results are shown in Table 2. As we can see from this table, when the number of registers increases, the Vysochanskij-Petunin and Chebyshev error bounds and the experimental error tighten rapidly, and the mean error decreases but keeps oscillating around 0, as expected from Theorem 1 of [10].

p	VP	Chebyshev	Mean error	Variance error
2^8	± 0.7506	± 1.367	$3.858 \cdot 10^{-3}$	0.01485
2^{10}	± 0.3219	± 0.5230	$3.987 \cdot 10^{-4}$	$3.236 \cdot 10^{-3}$
2^{12}	± 0.1520	± 0.2377	$2.945 \cdot 10^{-4}$	$7.631 \cdot 10^{-4}$
2^{14}	± 0.07560	± 0.1155	$-8.756 \cdot 10^{-6}$	$1.932 \cdot 10^{-4}$
2^{16}	± 0.03929	± 0.05954	$-6.172 \cdot 10^{-5}$	$4.950 \cdot 10^{-5}$
2^{18}	± 0.02158	± 0.03285	$2.175 \cdot 10^{-5}$	$1.183 \cdot 10^{-5}$

Table 2: Average error bounds and experimental bounds for conductance in the ball subgraph $\mathcal{B}_1(v)$ in 100 generated LFR-1 graphs with different numbers of registers.

5.1 Triangles

In Figure 3 we show the number of triangles and their estimates in the LFR-3 graph in ball subgraphs of radii 1, 2, and 3. Since the estimate for the number of triangles again has a low p -value on the DIP test for unimodality ($p < 0.01$), we can again use the Vysochanskij-Petunin bounds. There is a large difference in the number of triangles between the balls of radius 2 and 3, which can be explained by the fact that ball subgraphs of radius 3 contain a large fraction of the entire graph.

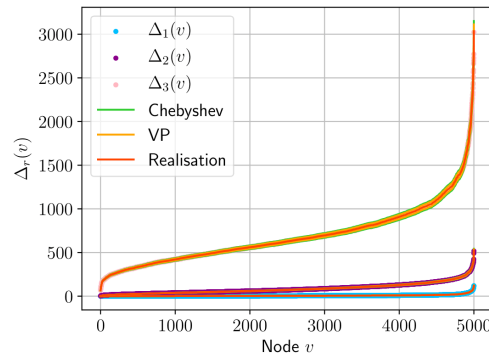


Figure 3: Estimate and realisation of the number of triangles $\Delta_r(v)$ in the ball subgraphs of radius $r = 1, 2, 3$ in the LFR-3 graph; the number of registers is $p = 2^{14}$. The nodes v are sorted by the realised number $\Delta_r(v)$ in ascending order.

5.1.1 Transitivity

In Figure 4, we show the transitivity in the ball subgraphs $\mathcal{B}_1(v)$ (Figure 4a) and $\mathcal{B}_2(v)$ (Figure 4b), together with the bounds from Theorem 8 and 9.

The DIP test for unimodality gives a p -value lower than 0.01 so the Vysochanskij-Petunin error bounds can again be used. Interestingly, the transitivity in ball subgraphs of radius 2 is already very small, which means that these ball subgraphs contain significantly more wedges than triangles. This is again an indication that the best communities in these kind of graphs are between the ball subgraphs of radius 1 and the ball subgraphs of radius 2, as also shown in [11]. Note that our estimation of the transitivity based on the HYPERBALL type algorithms is very precise, in fact the precision is much higher than suggested by the error bounds. This large precision may be explained by the fact that the number of wedges in a graph is often very large. This improves the precision of the transitivity estimator as well.

6 Application to community detection

In this section we will show how our HYPERBALL-based algorithms can help to detect communities in real-world networks. As an example, we use two large networks with a community structure: COM-DBLP and COM-AMAZON [16, 26]. Table 3 summarises the properties of these graphs.

We will identify the communities in these networks using the PAGERANK-NIBBLE algorithm [1] that detects a community by finding a set of nodes with low conductance, starting from a random seed set of

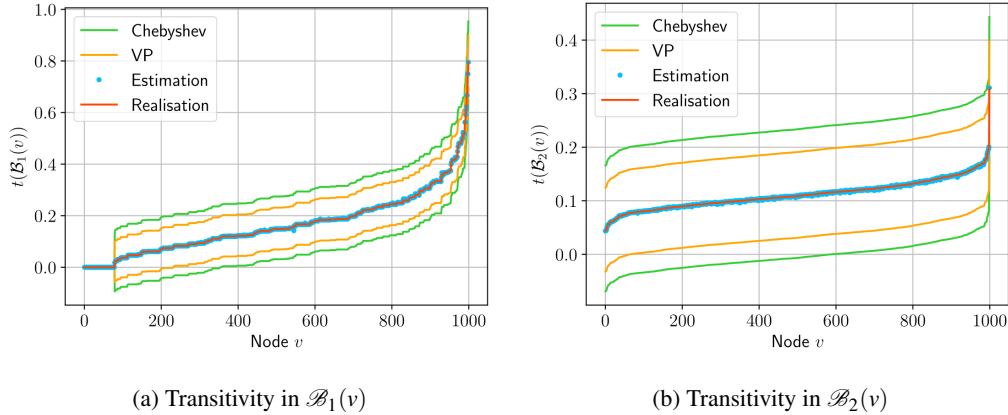


Figure 4: Transitivity estimate in a LFR-1 graph with $p = 2^{14}$ registers. The nodes v are sorted by realised transitivity in $\mathcal{B}_r(v)$, $r = 1, 2$, in ascending order.

	$n = V $	$ E $	$ \Delta $
COM-DBLP	317k	1M	2M
COM-AMAZON	334k	925k	667k

Table 3: Properties of the used real-world graphs

nodes. We propose to enhance this algorithm by using alternative seed sets found by our HYPERBALL algorithms. For example, a seed set can consist of ball subgraphs with the smallest conductance, or the largest density of triangles, or the largest transitivity. We have implemented this approach using 5 different seed sets, each of 100 nodes, in PAGERANK-NIBBLE:

1. ϕ -seeds: nodes v with smallest conductance in $\mathcal{B}_1(v)$ and $\mathcal{B}_2(v)$;
2. Δ -seeds: nodes v with largest number of triangles of radius 0 and radius 1, $|\Delta_0(v)|$ and $|\Delta_1(v)|$;
3. t -seeds: nodes v with highest transitivity in $\mathcal{B}_1(v)$ and $\mathcal{B}_2(v)$;
4. degree seeds: nodes with highest degree;
5. random seeds: randomly chosen nodes.

For the PAGERANK-NIBBLE algorithm, we have used a maximum cut size of 200, a teleport probability of $\alpha = 0.85$ and we calculate the ε -approximate PageRank vector with $\varepsilon = 10^{-8}$ [1]. We then compare the resulting conductance obtained with seed sets listed above.

The results are presented in Figures 5 and 6. For both graphs the lowest conductance subgraphs are found with the ϕ -seeds. We notice that the found sets of small conductance are often small nearly isolated sets of nodes. Interestingly, using the ϕ -seeds, PAGERANK-NIBBLE is able to find such sets, while with degree seeds it fails to do so.

In the Amazon graph, t -seeds result in low conductance subgraphs after PAGERANK-NIBBLE, suggesting that the local high transitivity is a good indication for a community. Notice that random seeds also yield low conductance sets, but t -seeds clearly outperform this benchmark.

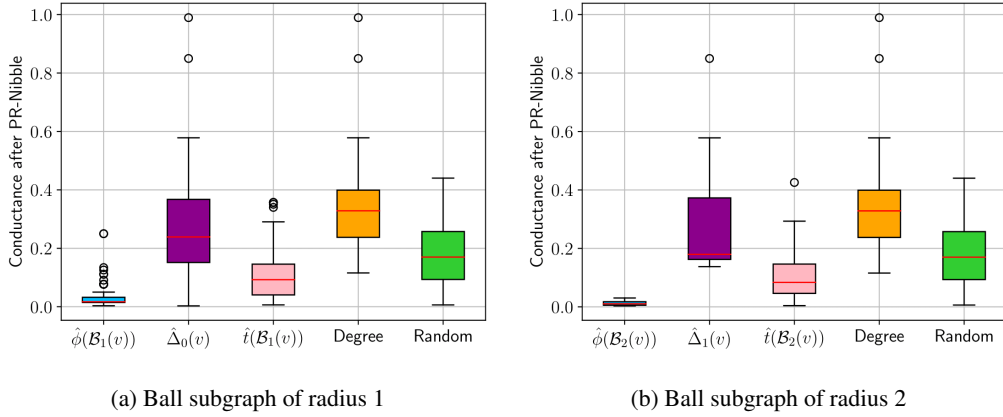


Figure 5: Boxplots of the resulting conductance after PR-Nibble in every seed set in the Amazon graph

In the DBLP graph, Figure 6, besides the ϕ -seeds, both Δ -seeds and t -seeds result in sets of low conductance after PAGERANK-NIBBLE. Moreover, when seeds are chosen based on ball subgraphs of radius 2, Δ -seeds work the best. Not only this helps us to detect communities, but we also obtain an insight that as communities have high number and density of triangles, it is important to use this knowledge to detect them. Figure 6 confirms that it is an important feature for community detection in this network: the t -seeds and the Δ -seeds result in a much lower conductance after PAGERANK-NIBBLE in comparison to the degree seeds and random seeds.

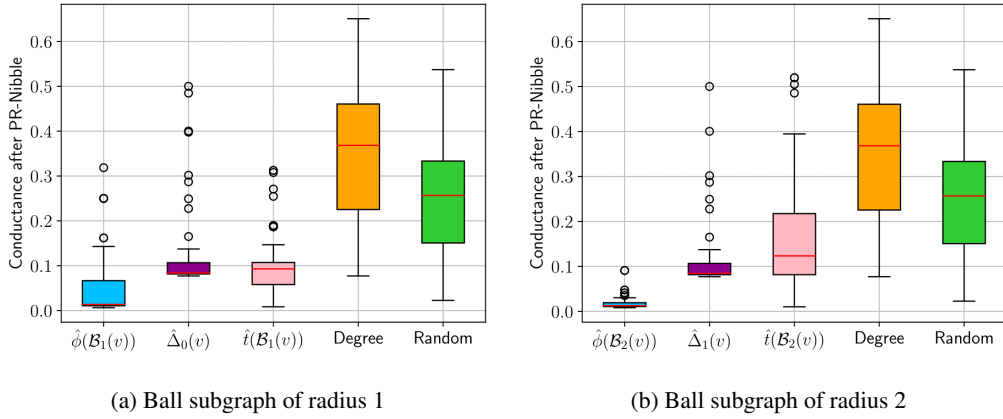


Figure 6: Boxplots of the resulting conductance after PR-Nibble in every seed set in the DBLP graph

We conclude that while the performance of different seed sets differ in different real-world networks, our HYPERBALL-based convincingly outperforms the random and degree-based seed sets.

7 Discussion

In this paper, we showed new applications of the HYPERBALL algorithm to count triangles and wedges. This enables us to approximate statistics like the transitivity and the conductance in ball subgraphs. Our estimates have good precision, for which we have derived explicit bounds. In this paper we demonstrated how these algorithms can be applied to choose good seed sets for community detection and to understand in more detail the structure of the communities.

Moreover, we showed that HYPERBALL can be extended to count the number of graphlets of any form in ball subgraphs. This gives us means to find the areas in large networks with high concentration of graphlets of specific kind. Potentially this will yield new ways to find anomalies in networks [9, 5], to distinguish the structure of networks of different nature (e.g. biological, technological or social networks) [17], and to compare real-world networks to mathematical random graph models, where the results on most likely locations of particular graphlets have been explicitly derived in recent literature [23, 19]. The bottleneck of the HYPERBALL-type algorithms for local graphlet count is the initialisation that requires, for each node v , the exact count of the graphlets that involve v . How to resolve this bottleneck remains an interesting question for further research.

References

- [1] Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using pagerank vectors. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pages 475–486. IEEE, 2006.
- [2] Lars Backstrom, Paolo Boldi, Marco Rosa, Johan Ugander, and Sebastiano Vigna. Four degrees of separation. In *Proceedings of the 4th Annual ACM Web Science Conference*, pages 33–42, 2012.
- [3] Vladimir Batagelj and Matjaž Zaveršnik. Generalized cores. *arXiv preprint cs/0202039*, 2002.
- [4] Luca Becchetti, Carlos Castillo, Debora Donato, Ricardo Baeza-Yates, and Stefano Leonardi. Link analysis for web spam detection. *ACM Transactions on the Web (TWEB)*, 2(1):1–42, 2008.
- [5] Luca Becchetti, Carlos Castillo, Debora Donato, Stefano Leonardi, and Ricardo A Baeza-Yates. Link-based characterization and detection of web spam. In *AIRWeb*, pages 1–8, 2006.
- [6] Yuchen Bian, Jingchao Ni, Wei Cheng, and Xiang Zhang. Many heads are better than one: Local community detection by the multi-walker chain. In *2017 IEEE International Conference on Data Mining (ICDM)*, pages 21–30. IEEE, 2017.
- [7] Paolo Boldi, Marco Rosa, and Sebastiano Vigna. Hyperanf: Approximating the neighbourhood function of very large graphs on a budget. In *Proceedings of the 20th international conference on World wide web*, pages 625–634, 2011.
- [8] Paolo Boldi and Sebastiano Vigna. In-core computation of geometric centralities with hyperball: A hundred billion nodes and beyond. In *2013 IEEE 13th International Conference on Data Mining Workshops*, pages 621–628. IEEE, 2013.
- [9] Lina Chen, Xiaoli Qu, Mushui Cao, Yanyan Zhou, Wan Li, Binhua Liang, Weiguo Li, Weiming He, Chenchen Feng, Xu Jia, et al. Identification of breast cancer patients based on human signaling network motifs. *Scientific reports*, 3:3368, 2013.

- [10] Philippe Flajolet, Eric Fusy, Olivier Gandouet, and Frédéric Meunier. Hyperloglog: the analysis of a near-optimal cardinality estimation algorithm. In *Analysis of Algorithms 2007 (AofA07)*, pages 127–146, 2007.
- [11] David F Gleich and C Seshadhri. Vertex neighborhoods, low conductance cuts, and good seeds for local community methods. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 597–605, 2012.
- [12] John A Hartigan, Pamela M Hartigan, et al. The dip test of unimodality. *The annals of Statistics*, 13(1):70–84, 1985.
- [13] JJ Huizinga. Estimating graph properties with hyperloglog-type algorithms. B.S. thesis, University of Twente, 2019.
- [14] Andrea Lancichinetti, Santo Fortunato, and Filippo Radicchi. Benchmark graphs for testing community detection algorithms. *Physical review E*, 78(4):046110, 2008.
- [15] Matthieu Latapy. Main-memory triangle computations for very large (sparse (power-law)) graphs. *Theoretical computer science*, 407(1-3):458–473, 2008.
- [16] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [17] Ron Milo, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824–827, 2002.
- [18] Satu Elisa Schaeffer. Graph clustering. *Computer science review*, 1(1):27–64, 2007.
- [19] Clara Stegehuis, Remco van der Hofstad, and Johan SH van Leeuwen. Variational principle for scale-free network motifs. *Scientific reports*, 9(1):1–10, 2019.
- [20] Mechthild Stoer and Frank Wagner. A simple min-cut algorithm. *Journal of the ACM (JACM)*, 44(4):585–591, 1997.
- [21] Ngoc Hieu Tran, Kwok Pui Choi, and Louxin Zhang. Counting motifs in the human interactome. *Nature communications*, 4(1):1–8, 2013.
- [22] Charalampos E Tsourakakis, Petros Drineas, Eirinaios Michelakis, Ioannis Koutis, and Christos Faloutsos. Spectral counting of triangles via element-wise sparsification and triangle-based link recommendation. *Social Network Analysis and Mining*, 1(2):75–81, 2011.
- [23] Remco van der Hofstad, Johan S. H. van Leeuwen, and Clara Stegehuis. Optimal subgraph structures in scale-free configuration models. *To appear in Annals of Applied Probability*, 2020.
- [24] DF Vysochanskij and Yu I Petunin. Justification of the 3σ rule for unimodal distributions. *Theory of Probability and Mathematical Statistics*, 21(25-36), 1980.
- [25] Pinghui Wang, Yiyang Qi, Yu Sun, Xiangliang Zhang, Jing Tao, and Xiaohong Guan. Approximately counting triangles in large graph streams including edge duplicates with a fixed memory usage. *Proceedings of the VLDB Endowment*, 11(2):162–175, 2017.
- [26] Jaewon Yang and Jure Leskovec. Defining and evaluating network communities based on ground-truth. *Knowledge and Information Systems*, 42(1):181–213, 2015.

A HyperLogLog and HyperBall algorithm

Algorithm A.1 The HYPERLOGLOG algorithm as described in [10], which approximates the cardinality of a data stream \mathcal{M} .

```
1: Let  $h_b(x)$  be the first  $b$  bits of the hashed value of element  $x$ 
2: Let  $h^b(x)$  be the other part of the hashed value of element  $x$ 
3: Let  $\rho(h^b(x))$  be the position of the leftmost 1-bit ( $\rho(001\dots) = 3$ ).
4:
5: initialise a collection of  $p = 2^b$  registers,  $M[1], \dots, M[p], to -\infty$ .
6:
7: function ADD( $M$ :counter,  $x$ : item)
8:    $i \leftarrow h_b(x)$ 
9:    $M[i] \leftarrow \max\{M[i], \rho(h^b(x))\}$ 
10: end function
11:
12: function SIZE( $M$ : counter)
13:    $Z \leftarrow (\sum_{j=0}^{p-1} 2^{-M[j]})^{-1}$ 
14:    $E \leftarrow \alpha_p p^2 Z$ 
15:   return  $E$ 
16: end function
17:
18: for each  $x \in \mathcal{M}$  do
19:   ADD( $M, x$ )
20: end for
21:
22: return SIZE( $M$ )
```

Algorithm A.2 The HyperBall algorithm as described in [8], which returns an estimation of the ball cardinality for each node. The functions ADD and SIZE of Algorithm A.1 are used.

```

1:  $c[-]$ , an array of  $n$  HyperLogLog counters initialised with nodes.
2:
3: function UNION( $M$ : counter,  $N$ : counter)
4:   for each  $i < p$  do
5:      $M[i] \leftarrow \max\{M[i], N[i]\}$ 
6:   end for
7: end function
8:
9: function COUNTBALL( $c$ : counter)
10:   $r \leftarrow 0$ 
11:  repeat
12:    for each  $v \in V$  do
13:       $a \leftarrow c[v]$ 
14:      for each  $w \in \mathcal{N}(v)$  do
15:         $a \rightarrow \text{UNION}(c[w], a)$ 
16:      end for
17:      write  $\langle v, a \rangle$  to disk
18:    end for
19:    update the array  $c[-]$  with the new  $\langle v, a \rangle$  pairs
20:     $r \leftarrow r + 1$ 
21:  until no counter changes its value
22:  return SIZE( $c$ )
23: end function
24:
25: for each  $v \in V$  do ▷ Initialisation
26:   ADD( $c[v], v$ )
27: end for
28:
29:  $(|\widehat{\mathcal{B}}_r|)_{r \geq 1} = \text{COUNTBALL}(c)$ 

```

B Proofs

B.1 Proof of Theorem 3

To find the error bound of this estimator $\hat{\phi}(\mathcal{B}_r(v))$ we use Chebyshev's inequality. The estimator $\hat{\phi}(\mathcal{B}_r(v))$ is a fraction of two other estimator, $|\widehat{\mathcal{E}}_r^+(v)|$ and $|\widehat{\mathcal{E}}_r^-(v)|$. The expectation and the variance of these estimators,

as the cardinality of the estimated set goes to infinity, are given in Theorem 1:

$$\mathbb{E}[\widehat{|\mathcal{E}_r(v)|}] \leq |\mathcal{E}_r(v)| \cdot (1 + \delta_1 + o(1)), \quad (23)$$

$$\mathbb{E}[\widehat{|\mathcal{E}_r^-(v)|}] \leq |\mathcal{E}_r^-(v)| \cdot (1 + \delta_1 + o(1)), \quad (24)$$

$$\text{Var}[\widehat{|\mathcal{E}_r(v)|}] \leq \eta^2 \cdot |\mathcal{E}_r(v)|^2, \quad (25)$$

$$\text{Var}[\widehat{|\mathcal{E}_r^-(v)|}] \leq \eta^2 \cdot |\mathcal{E}_r(v)|^2. \quad (26)$$

The coefficient η is defined as follows:

$$\eta = \frac{\beta_p}{\sqrt{p}} + \delta_2 + o(1). \quad (27)$$

The upper bounds are derived by using the fact that $|\delta_1(x)| \leq 5 \cdot 10^{-5} = \delta_1$ for all x and $|\delta_2(x)| \leq 5 \cdot 10^{-4} = \delta_2$ for all x , when the number of registers is larger or equal to 2^4 (Theorem 1). Throughout this work we assume that the number of registers is always larger than 2^4 .

Our goal is to use these estimators to find a lower and upper bound for $\hat{\phi}(\mathcal{B}_r(v))$. Following Chebyshev's theorem, we get the following inequalities for $p_1, p_2 > 0$ when the number of edges and directed edges in $\mathcal{B}_r(v)$ tend to infinity:

$$P\left(\left|\widehat{|\mathcal{E}_r(v)|} - \mathbb{E}[\widehat{|\mathcal{E}_r(v)|}]\right| \geq p_1\right) \leq \frac{|\mathcal{E}_r(v)|^2 \cdot \eta^2}{p_1^2}, \quad (28)$$

$$P\left(\left|\widehat{|\mathcal{E}_r^-(v)|} - \mathbb{E}[\widehat{|\mathcal{E}_r^-(v)|}]\right| \geq p_2\right) \leq \frac{|\mathcal{E}_r^-(v)|^2 \cdot \eta^2}{p_2^2}, \quad (29)$$

with η as defined in (27). We can rewrite the left-hand side of (28) using (23):

$$\begin{aligned} P\left(\left|\widehat{|\mathcal{E}_r(v)|} - \mathbb{E}[\widehat{|\mathcal{E}_r(v)|}]\right| \geq p_1\right) &= P\left(\left|\widehat{|\mathcal{E}_r(v)|} - |\mathcal{E}_r(v)| \cdot (1 + \delta_1(|\mathcal{E}_r(v)|) + o(1))\right| \geq p_1\right) \\ &= P\left(\left|\frac{\widehat{|\mathcal{E}_r(v)|} - |\mathcal{E}_r(v)| \cdot (1 + \delta_1(|\mathcal{E}_r(v)|) + o(1))}{|\mathcal{E}_r(v)|}\right| \geq \frac{p_1}{|\mathcal{E}_r(v)|}\right) \\ &= P\left(\frac{\widehat{|\mathcal{E}_r(v)|}}{|\mathcal{E}_r(v)|} \notin \left(1 + \delta_1(|\mathcal{E}_r(v)|) + o(1) - \frac{p_1}{|\mathcal{E}_r(v)|}, 1 + \delta_1(|\mathcal{E}_r(v)|) + o(1) + \frac{p_1}{|\mathcal{E}_r(v)|}\right)\right) \\ &\leq P\left(\frac{\widehat{|\mathcal{E}_r(v)|}}{|\mathcal{E}_r(v)|} \notin (1 - \varepsilon, 1 + \varepsilon)\right) \leq \frac{|\mathcal{E}_r(v)|^2 \cdot \eta^2}{p_1^2}, \end{aligned} \quad (30)$$

with $\varepsilon = \frac{p_1}{|\mathcal{E}_r(v)|} + \delta_1 + o(1)$. The first inequality follows because $\delta_1(x) \leq \delta_1$ for all x , and the last inequality directly follows from (28). Similarly, using (29), we obtain that for $\gamma = \frac{p_2}{|\mathcal{E}_r^-(v)|} + \delta_1 + o(1)$,

$$P\left(\frac{\widehat{|\mathcal{E}_r^-(v)|}}{|\mathcal{E}_r^-(v)|} \notin (1 - \gamma, 1 + \gamma)\right) \leq \frac{|\mathcal{E}_r^-(v)|^2 \cdot \eta^2}{p_2^2}, \quad (31)$$

In order to find an error bound for the conductance estimate (13) instead of the two estimators $|\widehat{\mathcal{E}_r(v)}|$ and $|\widehat{\mathcal{E}_r^-(v)}|$, we define the following two events:

$$A = \left\{ \frac{|\widehat{\mathcal{E}_r(v)}|}{|\mathcal{E}_r(v)|} \in (1 - \varepsilon, 1 + \varepsilon) \right\}, \quad (32)$$

$$B = \left\{ \frac{|\widehat{\mathcal{E}_r^-(v)}|}{|\mathcal{E}_r^-(v)|} \in (1 - \gamma, 1 + \gamma) \right\}. \quad (33)$$

By using the union bound and (30) and (31), we can express the intersection of events A and B as follows:

$$P(A \cap B) = 1 - P(\bar{A} \cup \bar{B}) \geq 1 - P(\bar{A}) - P(\bar{B}) \geq 1 - \frac{|\mathcal{E}_r(v)|^2 \cdot \eta^2}{p_1^2} - \frac{|\mathcal{E}_r^-(v)|^2 \cdot \eta^2}{p_2^2}. \quad (34)$$

Now, we rewrite the event $A \cap B$ to obtain probabilistic bounds for $\hat{\phi}(\mathcal{B}_r(v))$:

$$\begin{aligned} P(A \cap B) &\leq P \left[\frac{\frac{|\widehat{\mathcal{E}_r(v)}|}{|\mathcal{E}_r(v)|}}{\frac{|\widehat{\mathcal{E}_r^-(v)}|}{|\mathcal{E}_r^-(v)|}} \in \left(\frac{1 - \varepsilon}{1 + \gamma}, \frac{1 + \varepsilon}{1 - \gamma} \right) \right] \\ &= P \left[\frac{\frac{|\widehat{\mathcal{E}_r(v)}|}{|\mathcal{E}_r^-(v)|}}{\frac{|\mathcal{E}_r(v)|}{|\mathcal{E}_r^-(v)|}} \in \left(\frac{1 - \varepsilon}{1 + \gamma}, \frac{1 + \varepsilon}{1 - \gamma} \right) \right] \\ &= P \left[2 \frac{|\widehat{\mathcal{E}_r(v)}|}{|\mathcal{E}_r^-(v)|} - 1 \in \left(\frac{1 - \varepsilon}{1 + \gamma} \cdot \left(2 \frac{|\mathcal{E}_r(v)|}{|\mathcal{E}_r^-(v)|} - 1 \right), \frac{1 + \varepsilon}{1 - \gamma} \cdot \left(2 \frac{|\mathcal{E}_r(v)|}{|\mathcal{E}_r^-(v)|} - 1 \right) \right) \right] \\ &= P \left[\hat{\phi}(\mathcal{B}_r(v)) \in \left(\frac{1 - \varepsilon}{1 + \gamma} \cdot \phi(S_r(v)), \frac{1 + \varepsilon}{1 - \gamma} \cdot \phi(S_r(v)) \right) \right]. \end{aligned} \quad (35)$$

Combining (34) and (35) results in the error bounds for the conductance estimator:

$$P \left[\hat{\phi}(\mathcal{B}_r(v)) \in \left(\frac{1 - \varepsilon}{1 + \gamma} \cdot \phi(S_r(v)), \frac{1 + \varepsilon}{1 - \gamma} \cdot \phi(S_r(v)) \right) \right] \geq 1 - \eta^2 \left(\frac{|\mathcal{E}_r(v)|^2}{p_1^2} + \frac{|\mathcal{E}_r^-(v)|^2}{p_2^2} \right). \quad (36)$$

B.2 Proof of Theorem 6

Proof. The estimator of the triangle count of a graph comes directly from the Algorithm A.2. This means that we know the expectation and variance of this estimator when the number of triangles goes to infinity (Theorem 1):

$$\mathbb{E}[\hat{\Delta}_r(v)] = \Delta_r(v) \cdot \left(1 + \delta_1(\Delta_r(v)) + o(1) \right), \quad (37)$$

$$\text{Var}[\hat{\Delta}_r(v)] = \Delta_r(v)^2 \cdot \left(\frac{\beta_p}{\sqrt{p}} + \delta_2(\Delta_r(v)) + o(1) \right)^2, \quad (38)$$

for two oscillating functions $\delta_1(x)$ and $\delta_2(x)$ with $|\delta_1(x)| < 5 \cdot 10^{-5} = \delta_1$ and $|\delta_2(x)| < 5 \cdot 10^{-4} = \delta_2$ as soon as $p \geq 16$ registers and $\beta_p \approx 1.03896$ for $p > 128$ registers.

The proof of this theorem is straightforward: we use Chebyshev's inequality for the triangle estimator, where we substitute the expectation and variance from Equations (37) and (38):

$$P\left(|\hat{\Delta}_r(v) - E(\hat{\Delta}_r(v))| \geq a\right) \leq \frac{\text{Var}(\hat{\Delta}_r(v))}{a^2}, \quad (39)$$

$$P\left(|\hat{\Delta}_r(v) - E(\hat{\Delta}_r(v))| \geq a\right) \leq \frac{\Delta_r(v)^2 \cdot \left(\frac{\beta_p}{\sqrt{p}} + \delta_2(\Delta_r(v)) + o(1)\right)^2}{a^2}, \quad (40)$$

$$P\left(\hat{\Delta}_r(v) \in \left(E(\hat{\Delta}_r(v)) - a, E(\hat{\Delta}_r(v)) + a\right)\right) \geq 1 - \frac{\left(\frac{\beta_p}{\sqrt{p}} + \delta_2(\Delta_r(v)) + o(1)\right)^2 \Delta_r(v)^2}{a^2}, \quad (41)$$

for $a > 0$. Since $\eta = \frac{\beta_p}{\sqrt{p}} + \delta_2 + o(1) \geq \frac{\beta_p}{\sqrt{p}} + \delta_2(\Delta_r(v)) + o(1)$, as defined in (27), this concludes the proof. \square

Theorems 8 and 9 are proved in the same way as Theorems 3 and 5, since in both cases we need to find error bounds of a ratio of two estimators.