

The Arboricity Captures the Complexity of Sampling Edges

Talya Eden

Tel Aviv University, Tel Aviv, Israel
talyaa01@gmail.com

Dana Ron

Tel Aviv University, Tel Aviv, Israel
danaron@tau.ac.il

Will Rosenbaum

Max Planck Institute for Informatics, Saarbrücken, Germany
will.rosenbaum@gmail.com

Abstract

In this paper, we revisit the problem of sampling edges in an unknown graph $G = (V, E)$ from a distribution that is (pointwise) almost uniform over E . We consider the case where there is some a priori upper bound on the arboricity of G . Given query access to a graph G over n vertices and of average degree d and arboricity at most α , we design an algorithm that performs $O\left(\frac{\alpha}{d} \cdot \frac{\log^3 n}{\varepsilon}\right)$ queries in expectation and returns an edge in the graph such that every edge $e \in E$ is sampled with probability $(1 \pm \varepsilon)/m$. The algorithm performs two types of queries: degree queries and neighbor queries. We show that the upper bound is tight (up to poly-logarithmic factors and the dependence in ε), as $\Omega\left(\frac{\alpha}{d}\right)$ queries are necessary for the easier task of sampling edges from any distribution over E that is close to uniform in total variational distance. We also prove that even if G is a tree (i.e., $\alpha = 1$ so that $\frac{\alpha}{d} = \Theta(1)$), $\Omega\left(\frac{\log n}{\log \log n}\right)$ queries are necessary to sample an edge from any distribution that is pointwise close to uniform, thus establishing that a poly($\log n$) factor is necessary for constant α . Finally we show how our algorithm can be applied to obtain a new result on approximately counting subgraphs, based on the recent work of Assadi, Kapralov, and Khanna (ITCS, 2019).

2012 ACM Subject Classification Theory of computation \rightarrow Graph algorithms analysis; Theory of computation \rightarrow Approximation algorithms analysis; Theory of computation \rightarrow Streaming, sublinear and near linear time algorithms; Theory of computation \rightarrow Sketching and sampling

Keywords and phrases sampling, graph algorithms, arboricity, sublinear-time algorithms

Digital Object Identifier 10.4230/LIPIcs.ICALP.2019.52

Category Track A: Algorithms, Complexity and Games

Related Version A full version of the paper is available at <https://arxiv.org/abs/1902.08086>.

Funding *Talya Eden*: Supported by the Azrieli fellowship program for graduate students, by the Sephora Scholarship and by the Weinstein Graduate Studies Prize.

Dana Ron: Partially supported by the Israel Science Foundation grant No. 1146/18.

1 Introduction

Let $G = (V, E)$ be a graph over n vertices and m edges. We consider the problem of sampling an edge in G from a pointwise-almost-uniform distribution over E . That is, for each edge $e \in E$, the probability that e is returned is $(1 \pm \varepsilon)/m$, where ε is a given approximation parameter. An algorithm for performing this task has random access to the vertex set $V = \{1, \dots, n\}$ and can perform queries to G . The allowed queries are (1) *degree queries*



© Talya Eden, Dana Ron, and Will Rosenbaum;
licensed under Creative Commons License CC-BY

46th International Colloquium on Automata, Languages, and Programming (ICALP 2019).

Editors: Christel Baier, Ioannis Chatzigiannakis, Paola Flocchini, and Stefano Leonardi;

Article No. 52; pp. 52:1–52:14



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



denoted $\deg(v)$ (what is the degree, $d(v)$, of a given vertex v) and (2) *neighbor queries* denoted $\text{nbr}(v, i)$ (what is the i^{th} neighbor of v).¹ We refer to this model as *the uniform vertex sampling model*.

Sampling edges almost uniformly is a very basic sampling task. In particular it gives the power to sample vertices with probability approximately proportional to their degree, which is a useful primitive. Furthermore, there are sublinear algorithms that are known to work when given access to uniform edges (e.g., [1, 2]) and can be adapted to the case when the distribution over the edges is almost uniform (see Section 1.4 for details). An important observation is that in many cases it is crucial that the sampling distribution is pointwise-close to uniform rather than close with respect to the Total Variation Distance (henceforth TVD) – see the discussion in [15, Sec. 1.1].

Eden and Rosenbaum [15] recently showed that $\Theta^*(\sqrt{m}/d)$ queries are both sufficient and necessary for sampling edges almost uniformly, where $d = 2m/n$ denotes the average degree in the graph. (We use the notation O^* to suppress factors that are polylogarithmic in n and polynomial in $1/\varepsilon$.) The instances for which the task is difficult (i.e., for which $\Omega(\sqrt{m}/d)$ queries are necessary), are characterized by having very dense subgraphs, i.e., a subgraph with average degree $\Theta(\sqrt{m})$. Hence, a natural question is whether it is possible to achieve lower query complexity when some a priori bound on the density of subgraphs is known. A well studied measure for bounded density “everywhere” is the graph *arboricity* (see Definition 1.2 below). Indeed there are many natural families of graphs that have bounded arboricity such as graphs of bounded degree, bounded treewidth or bounded genus, planar graphs, graphs that exclude a fixed minor and many other graphs. In the context of social networks, preferential attachment graphs and additional generative models exhibit bounded arboricity [3, 6, 5], and this has also been empirically validated for many real-world graphs [18, 16, 22].

We describe a new algorithm for sampling edges almost uniformly whose runtime is $O^*(\alpha/d)$ where α is an upper bound on the arboricity of G . In the extremal case that $\alpha = \Theta(\sqrt{m})$, the runtime of our algorithm is the same as that of [15] (up to poly-log factors). For smaller α , our algorithm is strictly faster. In particular for $\alpha = O(1)$, the new algorithm is *exponentially* faster than that of [15]. We also prove matching lower bounds, showing that for all ranges of α , our algorithm is query-optimal, up to polylogarithmic factors and the dependence in $1/\varepsilon$.

Furthermore, while not as simple as the algorithm of [15], our algorithm is still easy to implement and does not incur any large constants in the query complexity and running time, thus making it suitable for practical applications.

1.1 Problem definition

In order to state our results precisely, we define the notion of pointwise-closeness of probability distributions (cf. [15]) and arboricity of a graph.

► **Definition 1.1.** Let D be a fixed probability distribution on a finite set X . We say that a probability distribution \hat{D} is **pointwise ε -close** to D if for all $x \in X$,

$$\left| \hat{D}(x) - D(x) \right| \leq \varepsilon D(x), \quad \text{or equivalently} \quad 1 - \varepsilon \leq \frac{\hat{D}(x)}{D(x)} \leq 1 + \varepsilon.$$

¹ If $i > d(v)$ then a special symbol, e.g. \perp , is returned.

If $D = U$, the uniform distribution on X , then we say that \widehat{D} is **pointwise ε -close to uniform**.

For the sake of conciseness, from this point on, unless explicitly stated otherwise, when we refer to an edge sampling algorithm, we mean an algorithm that returns edges according to a distribution that is pointwise-close to uniform.

► **Definition 1.2.** Let $G = (V, E)$ be an undirected graph. A forest $F = (V_F, E_F)$ (i.e., a graph containing no cycles) with vertex set $V_F = V$ and edge set $E_F \subseteq E$ is a **spanning forest** of G . We say that a family of spanning forests F_1, F_2, \dots, F_k **covers** G if $E = \bigcup_{i=1}^k E_{F_i}$. The **arboricity** of G , denoted $\alpha(G)$, is the minimum k such that there exists a family of spanning forests of size k that covers G .

An edge-sampling algorithm is given as input an approximation parameter $\varepsilon \in (0, 1)$ and a parameter α which is an upper bound on the arboricity of G . The algorithm is required to sample edges according a distribution that is pointwise ε -close to uniform. To this end the algorithm is given query access to G . In particular we consider the aforementioned uniform vertex sampling model.

1.2 Results

We prove almost matching upper and lower bounds on the query complexity of sampling an edge according to a distribution that is pointwise-close to uniform when an upper bound on the arboricity of the graph is known. The first lower bound stated below (Theorem 2) holds even for the easier task of sampling from a distribution that is close to uniform in TVD.

► **Theorem 1.** *There exists an algorithm \mathcal{A} that for any n, m, α , and graph $G = (V, E)$ with n nodes, m edges, and arboricity at most α , satisfies the following. Given n and α , \mathcal{A} returns an edge $e \in E$ sampled from a distribution \widehat{U} that is pointwise ε -close to uniform using $O\left(\frac{\alpha}{d} \cdot \frac{\log^3 n}{\varepsilon}\right)$ degree and neighbor queries in expectation, where $d = m/n$.*

In Section 1.5.1 we provide a high-level presentation of the algorithm referred to in Theorem 1 and discuss how it differs from the algorithm in [15] (for the case that there is no given upper bound on the arboricity).

We next state our lower bound, which matches the upper bound in Theorem 1 up to a polylogarithmic dependence on n (for constant ε).

► **Theorem 2.** *Fix $\varepsilon \leq 1/6$ and let n, m and α be parameters such that $\alpha = \alpha(n) \leq \sqrt{m}$ and $m \leq n\alpha$. Let $\mathcal{G}_{n,m}^\alpha$ be the family of graphs with n vertices, m edges and arboricity at most α . Then any algorithm \mathcal{A} that for any $G \in \mathcal{G}_{n,m}^\alpha$ samples edges in G from a distribution that is ε -close to uniform in total variation distance – and in particular, any distribution that is pointwise ε -close to uniform – requires $\Omega(\alpha/d)$ queries in expectation.*

When α is a constant then (assuming that $m = \Omega(n)$) the lower bound in Theorem 2 is simply $\Omega(1)$, while Theorem 1 gives an upper bound of $O(\log^3 n)$ (for constant ε). We prove that an almost linear dependence on $\log n$ is necessary, even for the case that $\alpha = 1$ (where the graph is a tree).

► **Theorem 3.** *Fix $\varepsilon \leq 1/6$, and let \mathcal{T}_n be the family of trees on n vertices. Then any algorithm \mathcal{A} that for any $G \in \mathcal{T}_n$ samples edges in G from a distribution that is pointwise ε -close to uniform requires $\Omega\left(\frac{\log n}{\log \log n}\right)$ queries in expectation.*

We note that both of our lower bounds hold when the algorithm is also given access to *pair queries*, $\text{pair}(u, v)$, which state whether u and v share an edge.

1.3 Discussion of the results

Comparison to known results. The simplest algorithm for sampling edges uniformly is based on rejection sampling. Namely, it repeats the following until an edge is output: Sample a uniform vertex u , flip a coin with bias $d(u)/d_{\max}$, where d_{\max} is the maximum degree, and if the outcome is HEADS, then output a random edge (u, v) incident to u . The expected complexity of rejection sampling grows like d_{\max}/d , that is, linearly with the maximum degree. As noted earlier, Eden and Rosenbaum [11] show that this dependence on the maximum degree is not necessary (for approximate sampling), as $O^*(\sqrt{m}/d)$ queries and time always suffice, even when the maximum degree is not bounded (e.g., is $\Theta(n)$). However, if the maximum degree is bounded, and in particular if $d_{\max} = o(\sqrt{m})$, then rejection sampling has lower complexity than the [11] algorithm (e.g., in the case that the graph is close to regular, so that when $d_{\max} = O(d)$, we get complexity $O(1)$).

Since the arboricity of a graph is both upper bounded by d_{\max} and by \sqrt{m} , our algorithm can be viewed as “enjoying both worlds”. Furthermore, our results can be viewed as showing that the appropriate complexity measure for sampling edges is not the maximum degree but rather the maximum average degree (recall that the arboricity α measures the maximum density of any subgraph of G – for a precise statement, see Theorem 4).

Sampling vs. Estimating. As shown by Eden, Ron and Seshadhri [11], $\Theta^*(\alpha/d)$ is also the complexity of estimating the number of edges in a graph when given a bound α on the arboricity of the graph. Their algorithm improves on the previous known bound of $O^*(\sqrt{m}/d)$ by Feige [17] and Goldreich and Ron [19] (when the arboricity is $o(\sqrt{m})$). However, other than the complexity, our algorithm for sampling edges and the algorithm of [11] for estimating the number of edges do not share any similarities, in particular, as the result of [11] is allowed to “ignore” an ε -fraction of the graph edges.

Furthermore, while the complexity of sampling and of approximate counting of edges are the same (up to $\log n$ and $1/\varepsilon$ dependencies), we have preliminary results showing that for other subgraphs this is not necessarily the case. Specifically, there exist graphs with constant arboricity for which estimating the number of triangles can be done using $O^*(1)$ queries, but pointwise-close to uniform sampling requires $\Omega(n^{1/4})$ queries in expectation.

On the necessity of being provided with an upper bound on the arboricity. While our algorithm does not require to be given any bound on the average degree d , it must be provided with an upper bound α on the arboricity of the given graph. To see why this is true, consider the following two graphs. The first graph consists of a perfect matching between its vertices, so that both its average degree and its arboricity are 1. For $\tilde{\alpha} > 1$, the second graph consists of a perfect matching over $n - n/\tilde{\alpha}$ vertices and an $\tilde{\alpha}$ -regular graph over the remaining $n/\tilde{\alpha}$ vertices. This graph has an average degree of roughly 2, and arboricity $\tilde{\alpha}$. If an edge-sampling algorithm is not provided with an appropriate upper bound on the arboricity, but is still required to run in (expected) time that grows like the ratio between the arboricity and the average degree, then the algorithm can be used to distinguish between the two graphs. However, assuming a random labeling of the vertices of the two graphs, this cannot be done in time $o(\tilde{\alpha})$.²

Pointwise closeness vs. closeness with respect to the TVD. The lower bound of Theorem 2 holds for sampling from a distribution that is close to uniform with respect to TVD,

² We note that this construction can be extended to work for any $d \leq \tilde{\alpha}$.

and *a fortiori* to sampling from pointwise-almost-uniform distributions. In contrast, the lower bound of Theorem 3 does not apply to sampling edges from a distribution that is ε -close to uniform in TVD. Indeed, a simple rejection sampling procedure (essentially ignoring all nodes with degrees greater than $1/\varepsilon$) can sample edges from a distribution that is ε -close to uniform in TVD using $O(1/\varepsilon)$ queries in expectation. Thus, Theorem 3 gives a separation between the tasks of sampling from distributions that are pointwise-close to uniform versus close to uniform in TVD. The general upper and lower bounds of Theorems 1 and 2 show that the separation between the complexity of these tasks can be at most poly-logarithmic for any graph.

1.4 An application to approximately counting subgraphs

In a recent paper [2], Assadi, Kapralov, and Khanna made significant progress on the question of counting arbitrary subgraphs in a graph in sublinear time. Specifically, they provide an algorithm that estimates the number of occurrences of any arbitrary subgraph H in G , denoted by $\#H$, to within a $(1 \pm \varepsilon)$ -approximation with high probability. The running time of their algorithm is $O^*\left(\frac{m^{\rho(H)}}{\#H}\right)$, where $\rho(H)$ is the fractional edge cover of H .³ Their algorithm assumes access to uniform edge samples in addition to degree, neighbor and pair queries. As noted in [2], their algorithm can be adapted to work with edge samples that are pointwise ε -close to uniform (where this is not true for edge samples that are only ε -close to uniform in TVD – e.g., when all the occurrences of H are induced by an ε -fraction of the edges). Invoking the algorithm of [2], and replacing each edge sample with an invocation of `Sample-edge` results in the following corollary.

► **Corollary 1.** *Let G be a graph $G = (V, E)$ with n nodes, m edges, and arboricity at most α . There exists an algorithm that, given $n, \alpha, \varepsilon \in (0, 1)$, a subgraph H and query access to G , returns a $(1 \pm \varepsilon)$ approximation of the number of occurrences of H in G , denoted $\#H$. The expected query complexity and running time of the algorithm are*

$$O^*\left(\min\left\{m, \frac{n\alpha \cdot m^{\rho(H)-1}}{\#H}\right\}\right) \quad \text{and} \quad O^*\left(\frac{n\alpha \cdot m^{\rho(H)-1}}{\#H}\right),$$

respectively, where $\rho(H)$ denotes the fractional edge cover of H , and the allowed queries are degree, neighbor and pair queries.

Thus, by combining our result with [2], we extend the known results for approximately counting the number of subgraphs in a graph in the uniform vertex sampling model. Furthermore, for graphs in which $m = \Theta(n\alpha)$, we obtain the same query complexity and running time of [2] without the assumption that the algorithm has access to uniform edge samples.

1.5 A high-level presentation of the algorithm and lower bounds

1.5.1 The algorithm

While our results concern *undirected* graphs $G = (V, E)$, it will be helpful to view each edge $\{u, v\} \in E$ as a pair of *ordered* edges (u, v) and (v, u) . Sampling an edge (almost) uniformly is equivalent to sampling a vertex with probability (almost) proportional to its degree. Hence

³ The fractional edge cover of a graph $H = (V_H, E_H)$ is a mapping $\psi : E_H \rightarrow [0, 1]$ such that for each vertex $a \in V_H$, $\sum_{e \in E_H, a \in e} \psi(e) \geq 1$. The fractional edge-cover number $\rho(H)$ of H is the minimum value of $\sum_{e \in E_H} \psi(e)$ among all fractional edge covers ψ .

we focus on the latter task. A single iteration of the algorithm we describe either returns a vertex or outputs FAIL. We show that the probability that it outputs FAIL is not too large, and that conditioned on the algorithm returning a vertex, each vertex v is returned with probability proportional to its degree, up to a factor of $(1 \pm \varepsilon)$.

Our starting point is a structural decomposition result for graphs with bounded arboricity (Lemma 2.3). Our decomposition defines a partition of the graph's vertices into *levels* L_0, L_1, \dots, L_ℓ . For parameters θ and β , the first layer L_0 consists of all vertices with degree at most θ , and for $i > 0$, level L_i contains all vertices v that do not belong to previous levels L_0, \dots, L_{i-1} , but have at least $(1 - \beta)d(v)$ neighbors in these levels. We prove that for any graph with arboricity at most α , for $\theta = \Theta(\alpha \log n / \varepsilon)$ and $\beta = \Theta(\varepsilon / \log n)$, there exists such a partition into layers with $\ell = O(\log n)$ levels. We stress that the algorithm does not actually construct such a partition, but rather we use the partition in our analysis of the algorithm.⁴

In order to gain intuition about the algorithm and its analysis, suppose that all vertices in L_0 have degree exactly θ , and that all edges in the graph are between vertices in consecutive layers. Consider the following *random walk* algorithm. The algorithm first selects an index $j \in [0, \ell]$ uniformly at random. It then selects a vertex u_0 uniformly at random. If $u_0 \in L_0$, then it performs a random walk of length j starting from u_0 (otherwise it outputs FAIL). If the walk did not pass through any vertex in L_0 (with the exception of the starting vertex u_0), then the algorithm returns the final vertex reached.

First observe that for every $u \in L_0$, the probability that u is returned is $\frac{1}{\ell+1} \cdot \frac{1}{n}$ (the probability that the algorithm selected $j = 0$ and selected u as u_0). This equals $\frac{d(u)}{(\ell+1) \cdot \theta \cdot n}$ (by our assumption that $d(u) = \theta$ for every $u \in L_0$). Now consider a vertex $v \in L_1$. The probability that v is returned is at least $\frac{1}{\ell+1} \cdot \frac{(1-\beta)d(v)}{n} \cdot \frac{1}{\theta} = \frac{(1-\beta)d(v)}{(\ell+1) \cdot \theta \cdot n}$ (the probability that the algorithm selected $j = 1$, then selected one of v 's neighbors $u \in L_0$, and finally selected to take the edge between u and v). In general, our analysis shows that for every i and every $v \in L_i$, the probability that v is returned is at least $\frac{(1-\beta)^i d(v)}{(\ell+1) \cdot \theta \cdot n}$. On the other hand, we show that for every vertex v , the probability that v is returned is at most $\frac{d(v)}{(\ell+1) \cdot \theta \cdot n}$. By the choice of θ and β we get that each vertex v is returned with probability in the range $[(1 - \varepsilon)d(v)\rho(\varepsilon, n), d(v)\rho(\varepsilon, n)]$ for $\rho(\varepsilon, n) = \Theta(\varepsilon / (\alpha n \log^2 n))$. By repeating the aforementioned random-walk process until a vertex is returned – $\Theta\left(\frac{\alpha n}{m} \cdot \frac{\log^2 n}{\varepsilon}\right) = \Theta\left(\frac{\alpha}{d} \cdot \frac{\log^2 n}{\varepsilon}\right)$ times in expectation – we obtain a vertex that is sampled with probability proportional to its degree, up to $(1 \pm \varepsilon)$.

We circumvent the assumption that $d(u) = \theta$ for every $u \in L_0$ by rejection sampling: In the first step, if the algorithm samples $u_0 \in L_0$, then it continues with probability $d(u_0)/\theta$ and fails otherwise. The assumption that all edges are between consecutive levels is not necessary for the analysis described above to hold. The crucial element in the analysis is that for every vertex $v \notin L_i$, where $i > 0$, at least $(1 - \beta)$ of the neighbors of v belong to L_0, \dots, L_{i-1} . This allows us to apply the inductive argument for the lower bound on the probability that v is returned when we average over all choices of j (the number of steps in the random walk). For precise details of the algorithm and its analysis, see Section 2.

We briefly discuss the relation between our algorithm for bounded-arboricity graphs, which we denote by \mathcal{A}_{ba} and the algorithm presented in [15] (for the case that no upper bound is given on the arboricity), which we denote by \mathcal{A}_{ua} . The algorithm \mathcal{A}_{ua} can be

⁴ This decomposition is related to the forest decomposition of Barenboim and Elkin [4]. The main difference, which is essential for our analysis, is that the partition we define is based on the number of neighbors that a vertex has in lower levels *relative* to its degree, while in [4] the partition is based on the absolute number of neighbors in higher levels.

viewed as considering a partition of the graph vertices into just two layers according to a degree threshold of roughly $\sqrt{m/\varepsilon}$. It performs a random walk similarly to \mathcal{A}_{ba} , but where the walk has either length 0 or 1. This difference in the number of layers and the length of the walk, is not only quantitative. Rather, it allows \mathcal{A}_{ua} to determine to which layer a vertex belongs simply according to its degree. This is not possible in the case of \mathcal{A}_{ba} (with the exception of vertices in L_0). Nonetheless, despite the apparent “blindness” of \mathcal{A}_{ba} to the layers it traverses in the random walk, we can show the following: Choosing the length of the random walk uniformly at random and halting in case the walk returns to L_0 , ensures that each vertex is output with probability approximately proportional to its degree.

1.5.2 The lower bounds

In order to prove Theorem 2, we employ the method of [14] – which builds upon the paradigm introduced in [7] – based on communication complexity. The idea of the proof is to reduce from the two-party communication complexity problem of computing the disjointness function. The reduction is such that (1) any algorithm that samples edges from an almost-uniform distribution reveals the value of the disjointness function with sufficiently large probability, and (2) every allowable query can be simulated in the two-party communication setting using little communication.

As opposed to the proof of the lower bound for general α , in the case of $\alpha = 1$ we did not find a way to employ the communication-complexity method. Instead, we design a direct, albeit somewhat involved, proof from first-principles.

Specifically, in order to prove Theorem 3, we consider a complete tree in which each internal vertex has degree $\log n$ (so that its depth is $\Theta\left(\frac{\log n}{\log \log n}\right)$). We then consider the family of graphs that correspond to all possible labelings of such a tree. As noted in Section 1.5.1, sampling edges almost uniformly is equivalent to sampling vertices with probability approximately proportional to their degree. In particular, in our construction, the label of the root should be returned with probability approximately $\log n/n$. We show that any algorithm that succeeds in returning the label of the root of the tree with the required probability must perform $\Omega\left(\frac{\log n}{\log \log n}\right)$ queries.

To this end we define a *process* \mathcal{P} that interacts with any algorithm \mathcal{A} , answering \mathcal{A} 's queries while constructing a uniform random labeling of the vertices and edges in the tree. The vertices of the tree are assigned random labels in $[n]$, and for each vertex v in the tree, its incident edges are assigned random labels in $[d(v)]$. We say that \mathcal{A} *succeeds*, if after the interaction ends, \mathcal{A} outputs the label of the root of the tree, as assigned by \mathcal{P} .

Let $L = \frac{\log n}{C \cdot \log \log n}$ be the lower bound we would like to prove, where C is a sufficiently large constant (so that in particular, L is a (small) constant fraction of the depth of the tree). Intuitively, \mathcal{A} would like to “hit” a vertex at depth at most L and then “walk up the tree” to the root. There are two sources of uncertainty for \mathcal{A} . One is whether it actually hits a vertex at depth at most L , and the second is which edges should be taken to go up the tree. Our lower bound argument mainly exploits the second uncertainty, as we sketch next.

The process \mathcal{P} starts with an unlabeled tree (of the aforementioned structure), and assigns labels to its vertices and edges in the course of its interaction with \mathcal{A} . Recall that \mathcal{P} answers the queries of \mathcal{A} while constructing a uniform labeling. Therefore, whenever \mathcal{A} asks a query involving a *new* label (i.e., that has not yet appeared in its queries or answers to them), the vertex to which this label is assigned, should be uniformly selected among all vertices that are not yet labeled. We shall say that a vertex is *critical* if its depth is at most L . As long as no critical vertex is hit, \mathcal{A} cannot reach the root. This implies that if no critical vertex is

hit in the course of its queries, then the probability that \mathcal{A} succeeds is $O(1/n)$.

While the probability of hitting a critical vertex is relatively small, it is not sufficiently small to be deemed negligible. However, suppose that \mathcal{A} hits a critical vertex u at depth $\Delta < L$, which occurs with probability roughly $(\log n)^\Delta/n$. Then, conditioned on this event, each of the $(\log n)^\Delta$ edge-labeled paths from u is equally likely to lead to the root, and the labels of vertices on these paths are uniformly distributed, thus intuitively conveying no information regarding the “right path”.

A subtlety that arises when formalizing this argument is the following. Suppose that in addition to hitting a critical vertex u , \mathcal{A} hits another vertex, v , which is not necessarily critical, but is at distance less than L from u (and in particular has depth at most $2L$, which we refer to as *shallow*). Then, a path starting from v might meet a path starting from u , hence adding a conditioning that makes the above argument (regarding uniform labelings) imprecise. We address this issue by upper bounding the probability of such an event (i.e., of hitting both a critical vertex and a shallow vertex), and accounting for an event of this type as a success of \mathcal{A} .

1.6 Related work

Some of the works presented below were already mentioned earlier in the introduction, but are provided in this subsection for the sake of completeness.

The work most closely related to the present work is the recent paper of Eden and Rosenbaum [15]. In [15], the authors proved matching upper and lower bounds of $\Theta^*(\sqrt{m}/d)$ for the problem of sampling an edge from an almost-uniform distribution in an arbitrary graph using degree, neighbor, and pair queries.

The problem of sampling edges in a graph is closely related to the problem of estimating m , the number of edges in the graph. In [17], Feige proved an upper bound of $O^*(\sqrt{m}/d)$ for obtaining a $(2 + \varepsilon)$ -factor multiplicative approximation of m using only degree queries,⁵ and shows that it is not possible to go below a factor of 2 with a sublinear number of degree queries. In [19], Goldreich and Ron showed that $\Theta^*(\sqrt{m}/d)$ queries are necessary and sufficient to obtain a $(1 + \varepsilon)$ -factor approximation of m if neighbor queries are also allowed.

Several works prove matching upper and lower bounds on the query complexity of counting the number of triangles [8], cliques [13], and star-graphs of a given size [20] using degree, neighbor, and pair queries (when the latter are necessary). Eden, Ron, and Seshadhri devised algorithms for estimating the number of k -cliques [12] and moments of the degree distribution [11] whose runtimes are parameterized by the arboricity α of the input graph (assuming a suitable upper bound for α is given to the algorithm as input). These algorithms outperform the lower bounds of [13] and [20] (respectively) in the case where $\alpha \ll \sqrt{m}$. In [9], Eden, Levi, and Ron described an efficient algorithm for distinguishing graphs with arboricity at most α from those that are far from any graph with arboricity 3α .

Two recent works [1, 2] consider a query model that allows uniform random edge sampling in addition to degree, neighbor, and pair queries. In this model, Aliakbarpour et al. [1] described an algorithm for estimating the number of star subgraphs. In the same model, Assadi et al. [2] devised an algorithm that relies on uniform edge samples as a basic query to approximately count the number of instances of an arbitrary subgraph in a graph. The

⁵ To be precise, Feige [17] shows that, given a lower bound d_0 on the average degree, $O(\sqrt{n/d_0}/\varepsilon)$ degree queries are sufficient. If such a lower bound is not provided to the algorithm, then a geometric search can be performed, as shown in [19].

results in [1] and [2] imply that uniform edge samples afford the model strictly more power: the sample complexity of the algorithm of [1] outperforms the lower bound of [20] for the same task, and the sample complexity of the algorithm of [2] outperforms the lower bound of [13] for estimating the number of cliques. (The results of [20, 13] are in the uniform vertex sampling model.)

Organization

Due to space constraints, in this extended abstract we provide full details only for the proof of our upper bound (Theorem 1). The proofs of Theorems 2 and 3 can be found in [10].

2 The Algorithm

In this section we describe an algorithm that samples an edge e from an arbitrary graph G with arboricity at most α , according to a pointwise almost uniform distribution. Theorem 1 follows from our analysis of the algorithm. In what follows, for integers $i \leq j$, we use $[i, j]$ to denote the set of integers $\{i, \dots, j\}$, and for a vertex v we let $\Gamma(v)$ denote its set of neighbors.

As noted in the introduction, sampling edges from a uniform distribution is equivalent to sampling vertices proportional to their degrees. Indeed, if each vertex v is sampled with probability $d(v)/2m$, then choosing a random neighbor $w \in \Gamma(v)$ uniformly at random returns the (directed) edge $e = (v, w)$ with probability $1/2m$. Thus, it suffices to sample each vertex $v \in V$ with probability (approximately) proportional to its degree.

2.1 Decomposing graphs of bounded arboricity

Before describing the algorithm, we describe a decomposition of a graph G into *layers* depending on its arboricity. We begin by recalling the following characterization of arboricity due to Nash-Williams [21].

► **Theorem 4** (Nash-Williams [21]). *Let $G = (V, E)$ be a graph. For a subgraph H of G , let n_H and m_H denote the number of vertices and edges, respectively, in H . Then $\alpha(G) = \max_H \{\lceil m_H / (n_H - 1) \rceil\}$, where the maximum is taken over all subgraphs H of G .*

► **Definition 2.1.** *Let $G = (V, E)$ be a graph, and $\theta \in \mathbf{N}$, $\beta \in (0, 1)$ parameters. We define a (θ, β) -**layering** in G to be the sequence of non-empty disjoint subsets $L_0, L_1, \dots, L_\ell \subseteq V$ defined by $L_0 = \{v \in V \mid d(v) \leq \theta\}$ and for $i \geq 1$,*

$$L_{i+1} = \{v \notin L_0 \cup L_1 \cup \dots \cup L_i \mid |\Gamma(v) \cap (L_0 \cup \dots \cup L_i)| \geq (1 - \beta)d(v)\}.$$

*That is, L_0 consists of all vertices of degree at most θ , and a vertex v is in L_{i+1} if i is the smallest index for which a $(1 - \beta)$ -fraction of v 's neighbors resides in $L_0 \cup L_1 \cup \dots \cup L_i$. We say that G admits a (θ, β) -**layered partition of depth ℓ** if we have $V = L_0 \cup L_1 \cup \dots \cup L_\ell$.*

► **Notation 2.2.** *For a fixed i , we denote $L_{\leq i} = L_0 \cup L_1 \cup \dots \cup L_i$, and similarly for $L_{< i}$, $L_{\geq i}$, and $L_{> i}$. We use the notation $d_i(v)$ to denote $|\Gamma(v) \cap L_i|$ and similarly for $d_{\leq i}(v)$ and $d_{\geq i}(v)$.*

► **Lemma 2.3.** *Suppose G is a graph with arboricity at most α . Then G admits a (θ, β) -layered partition of depth ℓ for $\theta = 4\alpha \lceil \log n \rceil / \varepsilon$, $\beta = \varepsilon / 2 \lceil \log n \rceil$, and $\ell \leq \lceil \log n \rceil$.*

Proof. For each i , let $W_i = V \setminus (L_0 \cup L_1 \cup \dots \cup L_{i-1})$ be the set of vertices not in levels $0, 1, \dots, i - 1$. Let $m(W_i)$ denote the number of edges in the subgraph of G induced by W_i .

52:10 The Arboricity Captures the Complexity of Sampling Edges

For any fixed i and $v \in W_{i+1}$, we have $d_{<i}(v) < (1 - \beta)d(v)$ because $v \notin L_{\leq i}$. Therefore, v has at least $\beta d(v) > \beta\theta$ neighbors in W_i . Summing over vertices $v \in W_{i+1}$ gives

$$m(W_i) = \frac{1}{2} \sum_{v \in W_i} d_{\geq i}(v) \geq \frac{1}{2} \sum_{v \in W_{i+1}} d_{\geq i}(v) > \frac{1}{2} |W_{i+1}| \cdot \beta\theta. \quad (1)$$

On the other hand, since G has arboricity at most α , Theorem 4 implies that $m(W_i) \leq \alpha |W_i|$. Plugging this into Equation (1), we find that $\frac{|W_{i+1}|}{|W_i|} \leq \frac{2\alpha}{\beta\theta} = \frac{1}{2}$, where the inequality is by the choice of β and θ . Therefore, for $\ell \leq \lceil \log n \rceil$, we have that $W_{\ell+1} = \emptyset$, implying that $V = L_0 \cup L_1 \cup \dots \cup L_\ell$. ◀

2.2 Algorithm description

The algorithm exploits the structure of graphs G with arboricity at most α described in Lemma 2.3. More precisely, as the algorithm does not have direct access to this structure, the structure is used explicitly only in the analysis of the algorithm. Let L_0, L_1, \dots, L_ℓ be a (θ, β) -layered partition of V with $\theta = 4\alpha \lceil \log n \rceil / \varepsilon$, $\beta = \varepsilon / \lceil \log n \rceil$, and $\ell = \lceil \log n \rceil$. Vertices $v \in L_0$ are sampled with probability exactly proportional to their degree using a simple rejection sampling procedure, **Sample- $L_0(G, \theta)$** . In order to sample vertices in layers L_i for $i > 0$, our algorithm performs a random walk starting from a random vertex in L_0 chosen with probability proportional to its degree. Specifically, the algorithm **Sample-edge(G, α)** chooses a length j to the random walk uniformly in $[0, \ell]$. The subroutine **Random-walk(G, θ, j)** performs the random walk for j steps, or until a vertex $v \in L_0$ is reached in some step $i > 0$. If the walk returns to L_0 , the subroutine aborts and does not return any vertex. (This behavior ensures that samples are not too biased towards vertices in lower layers.) Otherwise, **Random-walk** returns the vertex at which the random walk halts. Our analysis shows that the probability that the random walk terminates at any vertex $v \in V$ is approximately proportional to $d(v)$ (Corollary 2.7), although **Sample-edge** may fail to return any edge with significant probability. Finally, we repeat **Sample-edge** until it successfully returns a vertex.

Sample-edge(G, α, ε)

1. Let $\theta = 4\alpha \lceil \log n \rceil / \varepsilon$ and let $\ell = \lceil \log n \rceil$.
2. Choose a number $j \in [0, \ell]$ uniformly at random.
3. Invoke **Random-walk(G, θ, j)** and let v be the returned vertex if one was returned. Otherwise, **return FAIL**.
4. Sample a uniform neighbor w of v and **return** $e = (v, w)$.

Random-walk(G, θ, j)

1. Invoke **Sample- $L_0(\theta)$** and let v_0 be the returned vertex if one was returned. Otherwise, **return FAIL**.
2. For $i = 1$ to j do
 - a. Sample a random neighbor v_i of v_{i-1} .
 - b. If $v_i \in L_0$ then **return FAIL**.
3. **Return** v_j .

Sample- $L_0(G, \theta)$

1. Sample a vertex $u \in V$ uniformly at random and query for its degree.
2. If $d(u) > \theta$ **return FAIL**.
3. **Return** u with probability $\frac{d(u)}{\theta}$, and with probability $1 - \frac{d(u)}{\theta}$ **return FAIL**.

► **Definition 2.4.** We let $P_j[v]$ denote the probability that *Random-walk* returns v , when invoked with parameters G , θ and $j \in [0, \ell]$. We also let $P_{\leq j}[v] \stackrel{\text{def}}{=} \sum_{i=0}^j P_i[v]$ and similarly for $P_{\geq j}$.

► **Lemma 2.5.** Let ℓ be as set in Step 1 of *Sample-edge* and let $P_{\leq j}$ be as defined in Definition 2.4. For all $v \in V$, $P_{\leq \ell}[v] \leq \frac{d(v)}{n\theta}$.

Proof. We argue by induction on j that for any $j \in [0, \ell]$, $P_{\leq j}[v] \leq d(v)/n\theta$. For the case $j = 0$, it is immediate from the description of *Random-walk* and *Sample- L_0* that $P_0[v] = d(v)/(n\theta)$ if $v \in L_0$ and $P_0[v] = 0$ otherwise. Further, for $v \in L_0$, due to Step 2b, $P_i(v) = 0$ for all $i > 0$, so that the lemma holds for all $v \in L_0$. Now suppose that for all $v \in V$ we have $P_{\leq j-1}(v) \leq d(v)/n\theta$. Then for any fixed $v \notin L_0$ we compute

$$\begin{aligned} P_{\leq j}[v] &= \sum_{i=1}^j P_i[v] = \sum_{i=1}^j \sum_{u \in \Gamma(v)} P_{i-1}[u] \frac{1}{d(u)} = \sum_{u \in \Gamma(v)} \frac{1}{d(u)} \sum_{i=0}^{j-1} P_i[u] \\ &= \sum_{u \in \Gamma(v)} \frac{1}{d(u)} P_{\leq j-1}[u] \leq \sum_{u \in \Gamma(v)} \frac{1}{d(u)} \frac{d(u)}{n\theta} = \frac{d(v)}{n\theta}. \end{aligned}$$

The second equality holds by the definition of *Random-walk*, and the one before the last inequality holds by the inductive hypothesis. ◀

► **Lemma 2.6.** Let ℓ be as set in Step 1 of *Sample-edge*. For every $j \in [\ell]$, $v \in L_j$ and $k \in [j, \ell]$, we have $P_{\leq k}[v] \geq \frac{(1-\beta)^j d(v)}{n\theta}$.

Proof. We prove the claim by induction on j . For $j = 0$ and $k = 0$, by the description of *Random-walk* and *Sample- L_0* , for every $v \in L_0$,

$$P_0[v] = \frac{d(v)}{n\theta}. \quad (2)$$

For $j = 0$ and $0 < k \leq \ell$,

$$P_{\leq k}[v] = \sum_{i=0}^k P_i[v] = P_0[v] + \sum_{i=1}^k P_i[v] = \frac{d(v)}{n\theta}, \quad (3)$$

where the last equality is due to Step 2b in *Random-walk*.

For $j = 1$ and $1 \leq k \leq \ell$, for every $v \in L_1$, according to Step 2b in the procedure *Random-walk*, $P_0[v] = 0$. Also, for every $u \notin L_0$, $P_0[u] = 0$, since by Step 2 in *Sample- L_0* it always holds that v_0 is in L_0 . Therefore,

$$P_1[v] = \sum_{u \in \Gamma(v)} P_0[u] \cdot \frac{1}{d(u)} = \sum_{u \in \Gamma(v) \cap L_0} P_0[u] \cdot \frac{1}{d(u)} = \sum_{u \in \Gamma(v) \cap L_0} \frac{d(u)}{n\theta} \cdot \frac{1}{d(u)} = \frac{d_0(v)}{n\theta}, \quad (4)$$

where the second to last inequality is by Equation (2). By the definition of L_1 , for every $v \in L_1$, $d_0(v) \geq (1-\beta)d(v)$, and it follows that $P_{\leq k}[v] \geq P_1[v] \geq (1-\beta)d(v)/(n\theta)$.

We now assume that the claim holds for all $i \leq j-1$ and $k \in [i, \ell]$, and prove that it holds for j and for every $k \in [j, \ell]$. By the induction hypothesis and the definition of L_j , for any $v \in L_j$ we have

$$\begin{aligned} P_{\leq k}[v] &\geq P_{\leq j}[v] = \sum_{u \in \Gamma(v)} P_{\leq j-1}[u] \cdot \frac{1}{d(u)} \geq \sum_{i=0}^{j-1} \sum_{u \in \Gamma(v) \cap L_i} P_{\leq j-1}[u] \cdot \frac{1}{d(u)} \\ &\geq \sum_{i=0}^{j-1} \sum_{u \in \Gamma(v) \cap L_i} \frac{(1-\beta)^i d(u)}{n\theta} \cdot \frac{1}{d(u)} \geq \frac{(1-\beta)^{j-1} d_{\leq j-1}(v)}{n\theta} \geq \frac{(1-\beta)^j d(v)}{n\theta}. \end{aligned}$$

Hence, the claim holds for every $j \in [\ell]$ for every $k \in [j, \ell]$. ◀

52:12 The Arboricity Captures the Complexity of Sampling Edges

► **Corollary 2.7.** *For any graph G with arboricity at most α , the procedure **Sample-edge** when invoked with G , α and ε , returns each edge in the graph with probability in the range $\left[\frac{1-\varepsilon/2}{\rho}, \frac{1}{\rho}\right]$ for $\rho = n\theta(\ell + 1)$, $\theta = 4\alpha \lceil \log n \rceil / \varepsilon$ and $\ell = \lceil \log n \rceil$.*

Proof. Consider a specific edge $e^* = (v^*, w^*)$, and let i be the index such that $v^* \in L_i$. By the description of the procedure **Sample-edge**, the procedure **Random-walk** is invoked with an index j that is chosen uniformly in $[0, \ell]$. Hence, the probability that v^* is returned by **Random-walk** in Step 3 is $\Pr[v = v^*] = \frac{1}{\ell+1} \sum_{j=0}^{\ell} P_j[v] = \frac{1}{\ell+1} P_{\leq \ell}[v]$. By Lemma 2.5, $P_{\leq \ell}[v] \leq \frac{d(v)}{n\theta}$, and by Lemma 2.6, $P_{\leq \ell}[v^*] \geq \frac{(1-\beta)^\ell d(v^*)}{n\theta}$, where the probability is over the random coins of the procedures **Sample-edge** and **Random-walk**. Hence,

$$\Pr[v = v^*] \in [(1-\beta)^\ell, 1] \cdot \frac{d(v^*)}{n\theta(\ell+1)},$$

implying that for $\rho = n\theta(\ell+1)$,

$$\Pr[(v^*, w^*) \text{ is the returned edge}] \in [(1-\beta)^\ell, 1] \cdot \frac{1}{n\theta(\ell+1)} \in \left[\frac{1-\varepsilon/2}{\rho}, \frac{1}{\rho}\right], \quad (5)$$

where the last inequality is by the setting of $\beta = \varepsilon/2 \lceil \log n \rceil$. ◀

Proof of Theorem 1. Consider the algorithm that repeatedly calls **Sample-edge**(G, α) until an edge e is successfully returned. For a single invocation of **Sample-edge** and fixed edge e let A_e denote the event that **Sample-edge** returns e . By Corollary 2.7 we have that $\Pr[A_e] \geq (1-\varepsilon)/n\theta(\ell+1)$. Further, for any edge $e' \neq e$ the events A_e and $A_{e'}$ are disjoint, so we bound

$$\Pr[\text{Sample-edge returns an edge}] = \Pr\left[\bigcup_{e \in E} A_e\right] = \sum_{e \in E} \Pr[A_e] \geq \frac{(1-\varepsilon)m}{n\theta(\ell+1)}.$$

The expected number of iterations until **Sample-edge** succeeds is the reciprocal of this probability, so

$$\mathbf{E}[\# \text{ invocations until success}] \leq \frac{n\theta(\ell+1)}{(1-\varepsilon)m} = O\left(\frac{n\alpha}{m\varepsilon} \cdot \log^2 n\right).$$

Since each invocation of **Sample-edge** uses $O(\log n)$ queries, the expected number of queries before an edge is returned is $O\left(\frac{n\alpha}{\varepsilon m} \cdot \log^3 n\right)$.

Finally, when conditioned on a successful invocation of **Sample-edge**, Corollary 2.7 implies that for any $e, f \in E$ the probabilities p_e, p_f of returning e and f , respectively, satisfy

$$1 - \varepsilon/2 \leq \frac{p_e}{p_f} \leq \frac{1}{1 - \varepsilon/2} \leq 1 + \varepsilon.$$

Therefore, the induced distribution P over edges returned by a successful invocation of **Sample-edge** is pointwise ε -close to uniform, which gives the desired result. ◀

References

- 1 Maryam Aliakbarpour, Amartya Shankha Biswas, Themis Gouleakis, John Peebles, Ronitt Rubinfeld, and Anak Yodpinyanee. Sublinear-Time Algorithms for Counting Star Subgraphs via Edge Sampling. *Algorithmica*, pages 1–30, 2017. doi:10.1007/s00453-017-0287-3.
- 2 Sepehr Assadi, Michael Kapralov, and Sanjeev Khanna. A Simple Sublinear-Time Algorithm for Counting Arbitrary Subgraphs via Edge Sampling. In *ITCS*, volume 124 of *LIPICs*, pages 6:1–6:20. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2019.
- 3 Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- 4 Leonid Barenboim and Michael Elkin. Sublogarithmic distributed MIS algorithm for sparse graphs using Nash-Williams decomposition. *Distributed Computing*, 22(5-6):363–379, 2010. doi:10.1007/s00446-009-0088-2.
- 5 Reinhard Bauer, Marcus Krug, and Dorothea Wagner. Enumerating and generating labeled k -degenerate graphs. In *Proceedings of the Meeting on Algorithm Engineering & Experiments*, pages 90–98. Society for Industrial and Applied Mathematics, 2010.
- 6 Michael Baur, Marco Gaertler, Robert Görke, Marcus Krug, and Dorothea Wagner. Generating graphs with predefined k -core structure. In *Proceedings of the European Conference of Complex Systems*. Citeseer, 2007.
- 7 Eric Blais, Joshua Brody, and Kevin Matulef. Property testing lower bounds via communication complexity. *Computational Complexity*, 21(2):311–358, 2012.
- 8 Talya Eden, Amit Levi, Dana Ron, and C Seshadhri. Approximately counting triangles in sublinear time. In *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*, pages 614–633. IEEE, 2015.
- 9 Talya Eden, Reut Levi, and Dana Ron. Testing bounded arboricity. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2081–2092, 2018. doi:10.1137/1.9781611975031.136.
- 10 Talya Eden, Dana Ron, and Will Rosenbaum. The Arboricity Captures the Complexity of Sampling Edges, 2019. arXiv:1902.08086.
- 11 Talya Eden, Dana Ron, and C. Seshadhri. Sublinear Time Estimation of Degree Distribution Moments: The Degeneracy Connection. In Ioannis Chatzigiannakis, Piotr Indyk, Fabian Kuhn, and Anca Muscholl, editors, *44th International Colloquium on Automata, Languages, and Programming (ICALP 2017)*, volume 80 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 7:1–7:13, Dagstuhl, Germany, 2017. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. doi:10.4230/LIPIcs.ICALP.2017.7.
- 12 Talya Eden, Dana Ron, and C. Seshadhri. Faster sublinear approximations of k -cliques for low arboricity graphs. *CoRR*, abs/1811.04425, 2018. arXiv:1811.04425.
- 13 Talya Eden, Dana Ron, and C. Seshadhri. On approximating the number of k -cliques in sublinear time. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 722–734, 2018. doi:10.1145/3188745.3188810.
- 14 Talya Eden and Will Rosenbaum. Lower Bounds for Approximating Graph Parameters via Communication Complexity. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2018, August 20-22, 2018 - Princeton, NJ, USA*, pages 11:1–11:18, 2018. doi:10.4230/LIPIcs.APPROX-RANDOM.2018.11.
- 15 Talya Eden and Will Rosenbaum. On Sampling Edges Almost Uniformly. In Raimund Seidel, editor, *1st Symposium on Simplicity in Algorithms (SOSA 2018)*, volume 61 of *OpenAccess Series in Informatics (OASICs)*, pages 7:1–7:9, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. doi:10.4230/OASICs.SOSA.2018.7.
- 16 David Eppstein and Darren Strash. Listing all maximal cliques in large sparse real-world graphs. In *International Symposium on Experimental Algorithms*, pages 364–375. Springer, 2011.

52:14 The Arboricity Captures the Complexity of Sampling Edges

- 17 Uriel Feige. On sums of independent random variables with unbounded variance and estimating the average degree in a graph. *SIAM Journal on Computing*, 35(4):964–984, 2006.
- 18 Gaurav Goel and Jens Gustedt. Bounded arboricity to determine the local structure of sparse graphs. In *International Workshop on Graph-Theoretic Concepts in Computer Science*, pages 159–167. Springer, 2006.
- 19 Oded Goldreich and Dana Ron. Approximating average parameters of graphs. *Random Struct. Algorithms*, 32(4):473–493, 2008.
- 20 Mira Gonen, Dana Ron, and Yuval Shavitt. Counting stars and other small subgraphs in sublinear-time. *SIAM Journal on Discrete Mathematics*, 25(3):1365–1411, 2011.
- 21 C. St. JA. Nash-Williams. Edge-disjoint spanning trees of finite graphs. *Journal of the London Mathematical Society*, 1(1):445–450, 1961.
- 22 Kijung Shin, Tina Eliassi-Rad, and Christos Faloutsos. Patterns and anomalies in k -cores of real-world graphs with applications. *Knowledge and Information Systems*, 54(3):677–710, 2018.