

Optimizing Random Walk Based Statistical Estimation Over Graphs via Bootstrapping

Hong Xie¹, Pei Yi, Yongkun Li², and John C. S. Lui³, *Fellow, IEEE*

Abstract—Graphs are commonly used in various applications such as online social networks (OSNs), E-commerce systems and social recommender systems. Random walk sampling is often used to conduct statistical estimation over such graphs. This paper develops an algorithmic framework to reduce the mean square error of such statistical estimation. Our algorithmic framework is inspired by that the mean square error can be decomposed into a sum of the bias and variance of the estimator. More specifically, we apply the bootstrapping technique to design a bias reduction algorithm. Our bias reduction algorithm only utilizes a small number of “valid” sub-samples, which can reduce more bias of the estimator but may increase the variance of the estimator significantly. We use multiple parallel random walks to reduce this variance such that it can be reduced to arbitrarily small by deploying a sufficient number of random walks. We provide theoretical guarantees and computational complexity analysis of our proposed bias reduction algorithms. Our algorithmic framework enables one to attain different trade-offs between the sample complexity (i.e., number of parallel random walks) and the mean square error of the statistical estimation. Also, the proposed bias reduction algorithm is generic and can be applied to optimize a large class of random walk sampling algorithms. To demonstrate the versatility of the framework, we apply it to optimize the Metropolis random walk and simple random walk sampling. Extensive experiments on four public datasets confirm the effectiveness and computational efficiency of our proposed algorithmic framework under the mean square metric and beyond.

Index Terms—Random walk, bootstrapping, graph, statistical estimation

1 INTRODUCTION

STATISTICAL estimation over graphs is a fundamental task in graph analytic problems [1], [2], [3], [4], [5], [6], [7], [8]. A number of statistical estimation problems over graphs have been studied ranging from estimating simple statistics such as node degree distribution, vertex label distribution, size estimation [3], [4], [9], [10], [11], etc., to sophisticated statistics such as classification, ranking and regression [2], [12], [13], [14]. It is not an easy task to conduct statistical estimation over graphs. First, graphs in applications are usually large in scale. For example, the number of monthly active users of Facebook has reached over two billion [15]. Second, the whole graph is usually not accessible to third-party agents. In many OSNs, only APIs are available for third-party agents to access the graph data. Random walk sampling is a mainstream method to address this challenge [9], [15], [16], [17], [18], [19]. We use the following simplified

- Hong Xie and Pei Yi are with the Chongqing Key Laboratory of Software Theory and Technology, Chongqing University, Chongqing 400044, China. E-mail: {xiehong2018, 201914131045}@cqu.edu.cn.
- Yongkun Li is with the School of Computer Science and Technology, University of Science and Technology of China, Hefei, Anhui 230026, China. E-mail: ykli@ustc.edu.cn.
- John C. S. Lui is with the Department of Computer Science and Engineering, Chinese University of Hong Kong, Hong Kong. E-mail: cslui@cse.cuhk.edu.hk.

Manuscript received 5 Mar. 2021; revised 30 Sept. 2021; accepted 27 Oct. 2021. Date of publication 10 Nov. 2021; date of current version 3 Feb. 2023.

The work of Hong Xie was supported by the National Natural Science Foundation of China under Grant 61902042, Chongqing Talents: Exceptional Young Talents Project cstc2021ycjh-bgzxm0195. The work of John C.S. Lui was supported in part by GRF 14200420.

(Corresponding author: Hong Xie.)

Recommended for acceptance by B. He.

Digital Object Identifier no. 10.1109/TKDE.2021.3126906

example to illustrate random walk based statistical estimation over graphs.

Example 1. A company needs to make decisions regarding whether to do advertisements over a social network. The social network is characterized by a graph $\mathcal{G} \triangleq (\mathcal{V}, \mathcal{E}, x)$, where \mathcal{V} denotes the user set, $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ denotes the edge set and the function x prescribes an attribute $x(v) \in \{1, \dots, 10\}$ for vertex $v \in \mathcal{V}$. The attribute $x(v)$ quantifies the degree of proneness of user v on advertisements over social networks. The company wants to know the mean α and standard deviation σ of the degree of proneness over the whole user population:

$$\alpha = \sum_{v \in \mathcal{V}} x(v) / |\mathcal{V}|, \sigma = \sqrt{\sum_{v \in \mathcal{V}} (x(v) - \alpha)^2 / |\mathcal{V}|}.$$

Suppose we use the Metropolis random walk sampling algorithm [9] to get samples from the graph \mathcal{G} (details in Section 4). Suppose we get $L \in \mathbb{N}_+$ samples U_1, \dots, U_L , where $U_i \in \mathcal{V}$. Then one can estimate the mean and standard deviation as: $\hat{\alpha} = \sum_{i=1}^L x(U_i) / L$, $\hat{\sigma} = \sqrt{\sum_{i=1}^L (x(U_i) - \hat{\alpha})^2 / L}$.

Example 1 illustrates one typical characteristic of random walk based statistical estimation over graphs, i.e., only a finite number (usually small number) of samples are generated to conduct the estimation [8], [20]. Note that Example 1 can be mapped to a broad class of applications. In particular, the attribute function x in Example 1 can be mapped to node degree, vertex label, etc., leading to practical graph mining problems. One central problem is how to improve the estimation accuracy under this finite sample setting. A number of random walk algorithms were proposed to solve this problem [8], [9], [10], [15], [21]. In other words, these works improve estimation accuracy via getting “better” samples. This paper aims to improve estimation accuracy

from an orthogonal perspective, i.e., we apply bootstrapping techniques to exploit the property of the statistic in estimation to improve estimation accuracy.

Our framework is inspired by bootstrapping techniques and recent graph analytic systems that enable one to run millions of random walks in parallel on consumer-level personal computers [22], [23], [24], [25], [26]. The accuracy of the estimators in Example 1 can be characterized by bias and variance. To illustrate, consider the estimator $\hat{\sigma}$. The mean square error of $\hat{\sigma}$ is denoted by $\text{MSE}(\hat{\sigma})$

$$\text{MSE}(\hat{\sigma}) \triangleq \mathbb{E}[(\hat{\sigma} - \sigma)^2] = \text{Var}[\hat{\sigma}] + (\text{Bias}(\hat{\sigma}))^2, \quad (1)$$

where $\text{Var}[\hat{\sigma}]$ and $\text{Bias}(\hat{\sigma})$ are defined as: $\text{Var}[\hat{\sigma}] \triangleq \mathbb{E}[(\hat{\sigma} - \text{Mean}(\hat{\sigma}))^2]$, $\text{Bias}(\hat{\sigma}) \triangleq \text{Mean}(\hat{\sigma}) - \sigma$, with $\text{Mean}(\hat{\sigma}) \triangleq \mathbb{E}[\hat{\sigma}]$. This implies that one can reduce the estimation error by reducing the variance or bias. Note that these observations hold for many statistics beyond the mean and standard deviation and hence we consider general statistics. We apply the bootstrapping technique to reduce the bias. Unlike most previous bootstrapping techniques that reduce the bias under the constraint of not changing the variance or only allowing it to increase slightly [27], our bias reduction algorithm only utilizes a small number of “valid” sub-samples, which can reduce more bias of the estimator but may increase the variance of the estimator significantly. Then we use multiple parallel random walks to reduce the variance. One may argue that this rises an issue of scalability. Fortunately, it can be addressed by recent graph analytic systems [22], [23], [24], [25], [26], which enable running millions of random walks in parallel on consumer-level personal computers. In fact, the variance can be reduced by averaging, i.e., $\text{Var}[\hat{\sigma}_{\text{mr}}] = \text{Var}[\hat{\sigma}]/n$, where $\hat{\sigma}_{\text{mr}} = (\hat{\sigma}^{(1)} + \dots + \hat{\sigma}^{(n)})/n$ denotes the average of n estimations of σ from n parallel random walks with the same initial point. Note that the bias is unchanged, i.e., $\text{Bias}(\hat{\sigma}_{\text{mr}}) = \text{Bias}(\hat{\sigma})$. Namely, one can reduce the variance to arbitrarily small by deploying a sufficient number of random walks. We develop an algorithmic framework to reduce both bias and variance of the estimator. Our contributions are the following:

- We develop an algorithmic framework to reduce the mean square error of random walk based statistical estimation over graphs. Our algorithmic framework provides a novel combination of *random walk sampling* and *bootstrapping technique*, and it enables one to attain different trade-offs between sample complexity (i.e., number of parallel random walks) and mean square error of the statistical estimation.
- We apply the bootstrapping technique to design a bias reduction algorithm. Our bias reduction algorithm only utilizes a small number of “valid” sub-samples, which can reduce more bias of the estimator but may increase the variance of the estimator significantly. To overcome this problem, we use multiple parallel random walks to reduce this variance, and show that it can be reduced to arbitrarily small by deploying a sufficient number of random walks. Our bias reduction algorithm is generic and can be applied to a large class of random walk sampling

algorithms and statistical estimation problems. We provide theoretical guarantees and computational complexity analysis of our proposed bias reduction algorithms.

- To demonstrate the versatility of our framework, we apply it to optimize the Metropolis random walk sampling and simple random walk sampling. Experiment results on four public datasets show that our algorithmic framework can reduce the bias of both random walks without bias reduction (or with a classical bias reduction method) by as high as around 80% (or 60%). We also achieve similar reduction on mean square error using only 1000 parallel random walks and a larger number of random walks can lead to a larger reduction. Our proposed bias reduction algorithm only incurs a negligible extra computation. Our proposed algorithms still have superior performance over baselines under both the relative mean absolute error and relative mean cubic absolute error metric.

2 MODEL & PROBLEM FORMULATION

2.1 The Graph Model

We consider an undirected graph with a finite set $\mathcal{V} \triangleq \{1, \dots, V\}$ of vertices, where $V \in \mathbb{N}_+$. Each vertex v can be a user in an OSN, or an item in a social recommender system, etc. Let $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ denote the edge set. As the graph \mathcal{G} is undirected, then $(u, v) \in \mathcal{E}$ implies that $(v, u) \in \mathcal{E}$. For example, an edge $(u, v) \in \mathcal{E}$ can indicate the friendship between u and v in a social network. We focus on the case that the graph \mathcal{G} is connected. Let $d(v) \in \mathbb{N}_+$ denote the degree of vertex v , formally $d(v) \triangleq |\{u | (v, u) \in \mathcal{E}\}|$. Let $\mathcal{N}(v) \subseteq \mathcal{V}$ denote the neighbor set of vertex v , formally $\mathcal{N}(v) \triangleq \{u | (v, u) \in \mathcal{E}\}$. One can observe that $d(v) = |\mathcal{N}(v)|$.

We consider a real value attribute. In particular, each vertex v is associated with a value $x(v) \in \mathcal{X}$ indicating the attribute, where $\mathcal{X} \subseteq \mathbb{R}$ denotes the value set. For example, the value $x(v)$ can denote the gender of vertex v , then $\mathcal{X} = \{-1(\text{male}), 1(\text{female})\}$. The value $x(v)$ can also denote the degree of vertex v , then $\mathcal{X} = \{1, \dots, d_{\max}\}$, where $d_{\max} = \max_{v \in \mathcal{V}} d(v)$. For presentation simplicity, this paper considers the case that the cardinality of \mathcal{X} is finite, i.e., $|\mathcal{X}| < \infty$. Our results can be generalized to continuous case straightforwardly, by changing some summations into integrations. Let μ denote a distribution over the attribute value set \mathcal{X} , which summarizes the collective attribute value over the whole vertex set. Formally, we define μ as: $\mu(y) = \sum_{v \in \mathcal{V}} \frac{\mathbb{1}_{\{x(v)=y\}}}{|\mathcal{V}|}$, $\forall y \in \mathcal{X}$. Namely, $\mu(y)$ is the fraction of vertices with value $y \in \mathcal{X}$. To simplify notations, we denote the undirected graph as $\mathcal{G} = (\mathcal{V}, \mathcal{E}, x)$.

2.2 The Estimator Model

We consider a class of statistics $\theta \in \mathbb{R}$ over the graph \mathcal{G} such that it can be expressed as a function of the distribution μ , i.e., $\theta = T(\mu)$, where T denotes a mapping function: $T: \mu \mapsto \mathbb{R}$. For example, the mean of value can be modeled as

$$T(\mu) = \sum_{y \in \mathcal{X}} \mu(y)y. \quad (2)$$

The standard deviation of value can be modeled by T as

$$T(\mu) = \sqrt{\sum_{y \in \mathcal{X}} \mu(y)(y - \bar{y})^2}, \quad (3)$$

where $\bar{y} = \sum_{y \in \mathcal{X}} \mu(y)y$. The standard deviation can be generalized to $T(\mu) = (\sum_{y \in \mathcal{X}} \mu(y)|y - \bar{y}|^c)^{1/c}$, where $c \in \mathbb{R}_+$. We like to remark that many statistical estimation problems over the graph \mathcal{G} produce an estimator being a function of the distribution μ , e.g., maximum likelihood estimation, regression, etc. Namely, the statistic $\theta = T(\mu)$ can model a large class of statistical estimation problems over the graph.

2.3 Problem Formulation

We consider a large scale graph \mathcal{G} and one has to use random walk sampling to estimate the statistic θ . This setting is adopted in many previous works [8], [15], [28]. We aim to design an algorithm denoted by \mathbb{A} to estimate the statistic θ via samples generated by a random walk sampling algorithm running on the graph \mathcal{G} . We consider two metrics in the design of \mathbb{A} . The first one is the sample complexity defined as

$$\text{SC}(\mathbb{A}) \triangleq \# \text{ of samples required by the algorithm.}$$

The second one is the estimation error. We consider the mean square error defined as $\text{MSE}(\hat{\Theta}) \triangleq \mathbb{E}[(\hat{\Theta} - \theta)^2]$, where $\hat{\Theta}$ denotes an estimator of θ produced by the algorithm \mathbb{A} . The objective is to design an algorithm \mathbb{A} to estimate the statistic θ attaining different trade-offs between the above two metrics.

3 ALGORITHMIC FRAMEWORK

3.1 Design of the Algorithmic Framework

We first define a random walk sampling oracle, which supports *sampling query* and *estimating query*.

Definition 1. A random walk oracle denoted by *RWOracle* is defined as a function such that:

- for each query with initial point $U_0 \in \mathcal{V}$ and sample length $L \in \mathbb{N}_+$, it returns L samples denoted by $U \triangleq (U_1, \dots, U_L)$:

$$U = \text{RWOracle.Sampling}(U_0, L);$$

- for each query with a sequence of L samples U , it returns an estimation of the statistic θ :

$$\hat{\Theta}_{\text{rw}}(U) = \text{RWOracle.Statistics}(U).$$

We defer details of the random walk oracle *RWOracle* to Section 4. Here we focus on applying it to design our algorithmic framework. Note that for the linear statistic, we usually have $\hat{\Theta}_{\text{rw}}(U)$ being asymptotically unbiased, i.e., the bias $(\mathbb{E}[\hat{\Theta}_{\text{rw}}(U)] - \theta)$ goes to zero when L goes to infinity (more details in Section 4). Unfortunately, in practice, we can only collect a finite number of samples, under which the bias is not negligible. To reduce the bias, let us define an oracle to bootstrap the bias of the estimator $\hat{\Theta}_{\text{rw}}(U)$ first.

Definition 2. The bootstrapping oracle denoted by *BootBiasOracle* is defined as an algorithm such that for each query with a sample sequence U and the corresponding random walk oracle *RWOracle*, it returns an estimation for the bias of $\hat{\Theta}_{\text{rw}}(U)$

$$\Delta(\hat{\Theta}_{\text{rw}}(U)) = \text{BootBiasOracle}(U, \text{RWOracle}),$$

where $\Delta(\hat{\Theta}_{\text{rw}}(U))$ denotes the estimated bias of $\hat{\Theta}_{\text{rw}}(U)$.

We defer details of *BootBiasOracle* to Section 5. Here, let us focus on applying it to design our algorithmic framework.

To present our algorithmic framework, we need the following notations. Denote $M \in \mathbb{N}_+$ initial points as: $U_0 \triangleq (U_{0,1}, \dots, U_{0,M})$, where $U_{0,m} \in \mathcal{V}$. Denote $N_m \in \mathbb{N}_+$ as the number of parallel random walks associated with initial point $U_{0,m}$, where $m \in \{1, \dots, M\}$. For simplicity, denote $N \triangleq (N_1, \dots, N_M)$. Denote $L_m \in \mathbb{N}_+$ as the length of each parallel random walk associated with initial point $U_{0,m}$. For simplicity, denote $L \triangleq (L_1, \dots, L_M)$. Algorithm 1 outlines a parallel algorithmic framework to estimate the statistic θ over the graph \mathcal{G} . In step 2 of Algorithm 1, we run $\sum_{m=1}^M N_m$ random walks in parallel via the random walk oracle *RWOracle.Sampling*. These random walks are organized into M groups. Group m has N_m parallel random walks and each random walk within this group has the same initial point $U_{0,m}$. After we obtain samples from these $\sum_{m=1}^M N_m$ random walks, in step 3 we apply the bootstrapping oracle *BootBiasOracle* to estimate the bias of each random walk sequence. We then use the estimated bias to debias the random walk estimator by deducting it from the estimator. Finally, return the average of debiased estimators as an estimator of the statistic θ . The notation U is used with several different subscript forms, but none of them is overloaded, because the bold case U corresponds to a vector, while the normal case U corresponds to a scalar.

Algorithm 1. Algorithmic Framework

- 1: **Input:** U_0, N, L , Random walk oracle *RWOracle*, Bootstrapping oracle *BootBiasOracle*
- 2: **Parallel random walk sampling:** run $\sum_{m=1}^M N_m$ random walks in parallel to get samples:

$$U_{m,n} \leftarrow \text{RWOracle.Sampling}(U_{0,m}, L_m), \\ \forall m = 1, \dots, M, n = 1, \dots, N_m$$

- 3: Estimate the bias via bootstrapping:

$$\Delta_{m,n} \leftarrow \text{BootBiasOracle}(U_{m,n}, \text{RWOracle}), \\ \forall m = 1, \dots, M, n = 1, \dots, N_m$$

- 4: Compute the debiased estimators:

$$\hat{\Theta}_{m,n} \leftarrow \text{RWOracle.Estimate}(U_{m,n}) - \Delta_{m,n}, \\ \forall m = 1, \dots, M, n = 1, \dots, N_m$$

- 5: **Return:** $\hat{\Theta} \leftarrow \frac{1}{M} \sum_{m=1}^M \sum_{n=1}^{N_m} \frac{\hat{\Theta}_{m,n}}{N_m}$
-

Remark. The bias estimation step, i.e., step 3, of Algorithm 1 incurs extra computation. One may suggest to use this extra computation to generate more samples instead. We argue that this extra computation incurred by our proposed bias estimation algorithm, i.e., Algorithm 5, can be negligible compared to the total computation of generating samples. First, experiment results in Section 2.2 of our supplementary file, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org>.
Authorized licensed use limited to: Chinese University of Hong Kong. Downloaded on July 08, 2023 at 09:50:17 UTC from IEEE Xplore. Restrictions apply.

org/10.1109/TKDE.2021.3126906, show that Algorithm 1 implemented with Algorithm 5 for bias estimation has nearly the same running time as baseline random walk sampling algorithms without bias estimation. This implies that our bias estimation algorithm, i.e., Algorithm 5, only introduces a negligible extra computation. This negligible extra computation is achieved at the case that the whole graph is loaded into the main memory. When the graph cannot be loaded into the main memory, the running time of generating samples will be much larger as this incurs extra running time to load part of the graph into the main memory. Furthermore, the computational cost of generating one sample involves sampling from neighbors. When the number neighbors is not small, the computational cost of generating samples can be much larger than that of estimating the bias.

3.2 Analysis of the Algorithmic Framework

To analyze the performance of Algorithm 1, we first decompose the mean square error into follows: $\text{MSE}(\hat{\Theta}) = \text{Var}[\hat{\Theta}] + [\text{Bias}(\hat{\Theta})]^2$, where $\text{Var}[\hat{\Theta}]$ and $\text{Bias}(\hat{\Theta})$ are defined as the variance and bias of the estimator $\hat{\Theta}$

$$\text{Var}[\hat{\Theta}] \triangleq \mathbb{E}[(\hat{\Theta} - \text{Mean}(\hat{\Theta}))^2], \quad \text{Bias}(\hat{\Theta}) \triangleq \text{Mean}(\hat{\Theta}) - \theta,$$

with $\text{Mean}(\hat{\Theta}) \triangleq \mathbb{E}[\hat{\Theta}]$ denoting the mean of $\hat{\Theta}$.

Note that for a given group of random walks with the same initial point $U_{0,m}$, the $\hat{\Theta}_{m,n}$ across $n = 1, \dots, N_m$ are independent and identically distributed. Thus, we denote the variance and bias of $\hat{\Theta}_{m,n}$ as

$$\sigma_m^2 \triangleq \text{Var}[\hat{\Theta}_{m,n}], \quad \delta_m \triangleq \text{Bias}(\hat{\Theta}_{m,n}), \quad \forall n = 1, \dots, N_m.$$

Lemma 1. *The variance and bias of $\hat{\Theta}$ produced by Algorithm 1 can be derived as $\text{Var}[\hat{\Theta}] = \frac{1}{M^2} \sum_{m=1}^M \frac{\sigma_m^2}{N_m}$, $\text{Bias}(\hat{\Theta}) = \sum_{m=1}^M \frac{\delta_m}{M}$. Furthermore, $\lim_{\forall m, N_m \rightarrow \infty} \text{MSE}(\hat{\Theta}) = [\text{Bias}(\hat{\Theta})]^2$.*

Due to page limit, all proofs are in our supplementary file, available online. Lemma 1 states closed-form expressions for the variance and bias of the estimator $\hat{\Theta}$ outputted by Algorithm 1. One can reduce the variance of $\hat{\Theta}$ to zero by increasing the number of parallel random walks N_m at each initial point m to infinity. The bias of $\hat{\Theta}$ is the average of the bias of multiple random walks associated with M initial points. The following lemma derives analytical expression for the sample complexity of Algorithm 1.

Lemma 2. *The sample complexity of Algorithm 1 is $\text{SC}(\mathbb{A}_{\text{Alg01}}) = \sum_{m=1}^M N_m L_m$, where $\mathbb{A}_{\text{Alg01}}$ denotes Algorithm 1.*

Lemma 2 states the sample complexity for Algorithm 1. The sample complexity increases linearly with the sample length L_m associated with each initial point $U_{0,m}$, and also increases linearly with the number of parallel random walks in each initial point. In real-world applications, one usually does not want to have a long sample sequence L_m as it is time-consuming. However, one may want to use a larger number of parallel random walks to decrease the estimation error. Hence, the analytical expressions derived in Lemmas 1 and 2 enable one to select the appropriate sample complexity to trade it for improving estimation accuracy.

The computational complexity of Algorithm 1 involves $\sum_{m=1}^M N_m$ queries of the random walk oracle `RWOracle` and bootstrapping oracle `BootBiasOracle`. We next design algorithms to implement the the random walk oracle `RWOracle` and bootstrapping oracle `BootBiasOracle` as well as study their computational complexity.

4 THE SAMPLING ORACLE

In this section, we present two random walk oracles, which are based on the simple random walk and Metropolis random walk respectively. Most graph sampling algorithms are the variants of these two random walk algorithms, and they can also be applied to design the random walk oracle similarly (please refer to Section 9 for details).

Algorithm 2. `RWOracle` Based on Metropolis Random Walk

```

1: function RWOracle.Sampling  $U_0, L$ 
2:   for  $\ell = 1, \dots, \tilde{L}$  do ▷ Burn-in period
3:      $U_\ell \leftarrow v$  with prob.  $\mathbb{P}[v|U_{\ell-1}]$  derived in Equation (4)
4:   for  $\ell = \tilde{L} + 1, \dots, \tilde{L} + L$  do ▷ sampling period
5:      $U_\ell \leftarrow v$  with prob.  $\mathbb{P}[v|U_{\ell-1}]$  derived in Equation (4)
6:   return  $(U_{\tilde{L}+1}, \dots, U_{\tilde{L}+L})$ 
7: function RWOracle.Estimate  $U$ 
8:   Estimate the probability mass:

```

$$\hat{\mu}(y) \leftarrow \frac{\sum_{i=1}^{\text{length}(U)} \mathbb{1}_{\{x(U_i)=y\}}}{\text{length}(U)}, \quad \forall y \in \mathcal{X}$$

```

9:   return  $T(\hat{\mu})$ 

```

Metropolis Random Walk. The Metropolis random walk is an application of the Metropolis-Hastings algorithm [29] to graph sampling. Under the Metropolis random walk algorithm, in each step, the walker moves to a neighbor of the current vertex with certain probability derived as follows:

$$\mathbb{P}[U_{\ell+1}|U_\ell] = \begin{cases} \frac{1}{d(U_\ell)} \min\left(1, \frac{d(U_\ell)}{d(U_{\ell+1})}\right), & \text{if } U_{\ell+1} \in \mathcal{N}(U_\ell), \\ 1 - \sum_{v \in \mathcal{N}(U_\ell)} \frac{1}{d(U_\ell)} \min\left(1, \frac{d(U_\ell)}{d(v)}\right), & \text{if } U_{\ell+1} = U_\ell. \end{cases} \quad (4)$$

Note that the graph \mathcal{G} is connected. Thus, the Metropolis random walk has a stationary distribution. Let π denote the stationary distribution of the Metropolis random walk. It can be derived as $\pi(v) = 1/|\mathcal{V}|$, $\forall v \in \mathcal{V}$. Furthermore, for any function $f: \mathcal{V} \rightarrow \mathbb{R}$, it holds that $\lim_{L \rightarrow \infty} \sum_{\ell=0}^L \frac{f(U_\ell)}{L} = \sum_{v \in \mathcal{V}} \frac{f(v)}{|\mathcal{V}|}$. Namely, one can use the samples from the Metropolis random walk to construct asymptotic unbiased estimators (when the sample length goes to infinity). It is important to remember that when the sample length is finite, the bias is non-zero.

Based on the Metropolis random walk, Algorithm 2 outlines a random walk oracle. Consider the `RWOracle.Sampling`(U_0, L) function, the algorithm first simulates the Metropolis random walk by $\tilde{L} \in \mathbb{N}_+$ steps, for the purpose of making the random mix. This period is also called the burn-in period. Then the algorithm simulates another L steps and returns them as the samples. The function `RWOracle.Estimate`(U) first estimates the distribution μ

using the samples U via simple averaging. Then, it uses the estimated distribution to produce an estimator of the statistic θ . Note that $\text{length}(U)$ denotes the number of elements in vector U .

Via some direct arithmetic calculations, the computational complexity of `RWOracle.Sampling`(U_0, L) implemented in Algorithm 2 can be derived as $O((\tilde{L} + L) \max_{v \in \mathcal{V}} d(v))$. The factor $\max_{v \in \mathcal{V}} d(v)$ shows in the computational complexity due to that drawing a random sample from the neighbors of a node is linear in the number of neighbors. The computational complexity of `RWOracle.Estimate`(U) implemented in Algorithm 2 is $O(\text{length}(U))$.

Simple Random Walk. Under the simple random walk [30], the walker moves each neighbor of the current node with equal probability, formally

$$\mathbb{P}[U_{i+1}|U_i] = \begin{cases} \frac{1}{d(U_i)}, & \text{if } v \in \mathcal{N}(U_i), \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

The stationary distribution of the simple random walk can be expressed as $\pi(v) = d(v)/(2|\mathcal{E}|)$ [30]. Then, for any function $f: \mathcal{V} \rightarrow \mathbb{R}$, it holds that $\lim_{L \rightarrow \infty} \sum_{\ell=0}^L \frac{f(U_\ell)/d(U_\ell)}{\sum_{\ell=0}^L 1/d(U_\ell)} = \sum_{v \in \mathcal{V}} \frac{f(v)}{|\mathcal{V}|}$. Based on this observation, similar with Metropolis random walk, Algorithm 3 outlines a simple random walk based `RWOracle`.

Algorithm 3. Simple Random Walk Based `RWOracle`

```

1: function RWOracle.Sampling $U_0, L$ 
2:   for  $\ell = 1, \dots, \tilde{L}$  do ▷ burn in period
3:      $U_\ell \leftarrow v$  with prob.  $\mathbb{P}[v|U_{\ell-1}]$  derived in Equation (5)
4:   for  $\ell = \tilde{L} + 1, \dots, \tilde{L} + L$  do ▷ sampling period
5:      $U_\ell \leftarrow v$  with prob.  $\mathbb{P}[v|U_{\ell-1}]$  derived in Equation (5)
6:   return  $(U_{\tilde{L}+1}, \dots, U_{\tilde{L}+L})$ 
7: function RWOracle.Estimate $U$ 
8:   Estimate the probability mass:

```

$$\hat{\mu}(y) \leftarrow \frac{\sum_{i=1}^{\text{length}(U)} \mathbb{1}_{\{x(U_i)=y\}}/d(U_i)}{\sum_{i=1}^{\text{length}(U)} 1/d(U_i)}, \forall y \in \mathcal{X}$$

```

9:   return  $T(\hat{\mu})$ 

```

With a similar analysis as Algorithm 2, one can obtain that the function `RWOracle.Sampling`(U_0, L) and `RWOracle.Estimate`(U) implemented in Algorithm 3 have the same computational complexity as that implemented in Algorithm 2. Note that one can easily extend the above random walk sampling oracles to other random walk algorithms, such as its sophisticated variants [15].

5 BOOTSTRAPPING BIAS

We first design an oracle to bootstrap the bias via Jackknife. More importantly, we derive sufficient conditions, under which this oracle can reduce the bias of estimator. These conditions reveal a restriction that this oracle may perform poorly when the graph is incomplete. Then we design another oracle to relieve this restriction, and theoretically show that it works even over incomplete graphs.

Authorized licensed use limited to: Chinese University of Hong Kong. Downloaded on July 08, 2023 at 09:50:17 UTC from IEEE Xplore. Restrictions apply.

5.1 Bootstrapping Bias via Jackknife

Recall the closed-form expression for estimator $\hat{\Theta}$, i.e., $\hat{\Theta} = \frac{1}{M} \sum_{m=1}^M \sum_{n=1}^{N_m} \frac{\Theta_{m,n}}{N_m}$ (i.e., line 5 in Algorithm 1). By the linearity of expectation, it boils down to estimate the bias for $\hat{\Theta}_{m,n}$. The $\hat{\Theta}_{m,n}$ is evaluated from the samples denoted by

$$U_{m,n} \triangleq \{U_{m,n}^{(1)}, \dots, U_{m,n}^{(L_m)}\}.$$

For simplicity of notation, let $U_{m,n}^{(-i)}$ denote a vector of samples excluding the i th sample $U_{m,n}^{(i)}$

$$U_{m,n}^{(-i)} \triangleq [U_{m,n}^{(1)}, \dots, U_{m,n}^{(i-1)}, U_{m,n}^{(i+1)}, \dots, U_{m,n}^{(L_m)}].$$

Namely, $U_{m,n}^{(-i)}$ is a $(L_m - 1)$ -sized sub-sample from $U_{m,n}$. In total, we have L_m such sub-samples: $U_{m,n}^{(-1)}, \dots, U_{m,n}^{(-L_m)}$.

From the $(L_m - 1)$ -sized sub-sample $U_{m,n}^{(-i)}$, we apply the random walk oracle to generate one estimator as

$$\hat{\Theta}_{m,n}^{(-i)} = \text{RWOracle.Estimate}(U_{m,n}^{(-i)}).$$

Applying the Jackknife, the estimate of the bias of $\hat{\Theta}_{m,n}$ is

$$\widehat{\text{Bias}}(\hat{\Theta}_{m,n}) = (L_m - 1) \left(\frac{\sum_{i=1}^{L_m} \hat{\Theta}_{m,n}^{(-i)}}{L_m} - \hat{\Theta}_{m,n} \right). \quad (6)$$

Note that $\widehat{\text{Bias}}(\hat{\Theta}_{m,n})$ estimates the bias fully relying on the samples $U_{m,n}$ itself. Based on $\widehat{\text{Bias}}(\hat{\Theta}_{m,n})$, we outline a bootstrapping oracle to estimate the bias in Algorithm 4.

Algorithm 4. `BootBiasOracle`($U_{m,n}, \text{RWOracle}$) via Jackknife

```

1: for  $i = 1, \dots, L$  do
2:    $U_{m,n}^{(-i)} \leftarrow [U_{m,n}^{(1)}, \dots, U_{m,n}^{(i-1)}, U_{m,n}^{(i+1)}, \dots, U_{m,n}^{(L_m)}]$ 
3:    $\hat{\Theta}_{m,n}^{(-i)} \leftarrow \text{RWOracle.Estimate}(U_{m,n}^{(-i)})$ 
4:    $\hat{\Theta}_{m,n} \leftarrow \text{RWOracle.Estimate}(U_{m,n})$ 
5: Estimate the bias

```

$$\widehat{\text{Bias}}(\hat{\Theta}_{m,n}) \leftarrow (L_m - 1) \left(\frac{\sum_{i=1}^{L_m} \hat{\Theta}_{m,n}^{(-i)}}{L_m} - \hat{\Theta}_{m,n} \right)$$

```

6: return  $\widehat{\text{Bias}}(\hat{\Theta}_{m,n})$ 

```

To illustrate, consider $L_m = 3$. Then, we have $U_{m,n} = \{U_{m,n}^{(1)}, U_{m,n}^{(2)}, U_{m,n}^{(3)}\}$. In total, there are three sub-samples, i.e., $U_{m,n}^{(-1)} = [U_{m,n}^{(2)}, U_{m,n}^{(3)}]$, $U_{m,n}^{(-2)} = [U_{m,n}^{(1)}, U_{m,n}^{(3)}]$ and $U_{m,n}^{(-3)} = [U_{m,n}^{(1)}, U_{m,n}^{(2)}]$. The estimator can be calculated as $\hat{\Theta}_{m,n}^{(-1)} = \text{RWOracle.Estimate}([U_{m,n}^{(2)}, U_{m,n}^{(3)}])$. Similarly, $\hat{\Theta}_{m,n} \leftarrow \text{RWOracle.Estimate}([U_{m,n}^{(1)}, U_{m,n}^{(2)}, U_{m,n}^{(3)}])$. The bias can be calculated as

$$\widehat{\text{Bias}}(\hat{\Theta}_{m,n}) = 2 \left(\frac{\hat{\Theta}_{m,n}^{(-1)} + \hat{\Theta}_{m,n}^{(-2)} + \hat{\Theta}_{m,n}^{(-3)}}{3} - \hat{\Theta}_{m,n} \right).$$

To analyze the theoretical guarantee of Algorithm 4, we next define a class of expandable T .

Definition 3. Let X_1, \dots, X_n denote n independent and identically distributed samples from the distribution μ . Let $T(\mu_n)$ denote an estimator of the statistic θ , where $\mu_n(y) = \sum_{i=1}^n \mathbb{1}_{\{X_i=y\}}/n$. The T is expandable if it satisfies $\mathbb{E}[T(\mu_n)] = \theta + \sum_{j=1}^{\infty} a_j(\mu)/n^j$, where $a_j(\mu) \in \mathbb{R}, \forall j = 1, \dots, \infty$ is independent of n .

Most statistics are expandable [31], i.e., mean, variance, most maximum likelihood estimators, etc. For example, when the T is the mean, i.e., derived in Equation (2), we have $\mathbb{E}[T(\mu_n)] = \theta$. When T is the variance, i.e., derived in Equation (3), we have $\mathbb{E}[T(\mu_n)] = \theta - \theta/n$. The following theorem states theoretical guarantee for Algorithm 4 under expandable statistic.

Theorem 1. *Under the `RWOracle` implemented in Algorithms 2 or 3. Suppose \mathcal{G} is complete and of large scale and the T is expandable with $a_1(\mu) > 0$. Algorithm 4 corrects the bias from $\text{Bias}(\hat{\Theta}_{m,n}) = O(1/L_m)$ to $\text{Bias}(\hat{\Theta}_{m,n}^{\text{JK}}) = O(1/L_m^2)$, where $\hat{\Theta}_{m,n}^{\text{JK}} \triangleq \hat{\Theta}_{m,n} - \widehat{\text{Bias}}(\hat{\Theta}_{m,n})$ denotes the corrected estimator.*

Theorem 1 states sufficient conditions under which Algorithm 4 corrects the bias of estimator $\hat{\Theta}_{m,n}$ from $O(1/L_m)$ to $O(1/L_m^2)$. One sufficient condition is that \mathcal{G} has to be a complete graph. However, in real-world applications, graphs are usually incomplete. In this case, we may not have theoretical guarantee for Algorithm 4. In Section 5.2, we design a variant of the Jackknife algorithm to address this limitation. The following theorem states the computational complexity of Algorithm 4.

Theorem 2. *Under the `RWOracle` implemented in Algorithms 2 or 3. The computational complexity of the bootstrapping oracle `BootBiasOracle`($U_{m,n}$, `RWOracle`) implemented in Algorithm 4 is $O(L_m^2)$.*

Theorem 2 states that the computational complexity of `BootBiasOracle`($U_{m,n}$, `RWOracle`) implemented in Algorithm 4 is quadratic in the number of samples in $U_{m,n}$.

5.2 Improve Accuracy via Sub-Sample Selection

One restriction of Algorithm 4 is that some of the $(L_m - 1)$ -sized sub-samples of $U_{m,n}$ are not valid random walk sequences when the graph is incomplete. When the graph is incomplete, it may happen that $U_{m,n}^{(-i)}$ is not a sample sequence generated by the random walk sampling algorithm. In particular, according to the transition probability in Equation (4), the walker can not move from vertex $U_{m,n}^{(i-1)}$ to $U_{m,n}^{(i+1)}$. Formally, we define valid and invalid sub-sample in the following definition.

Definition 4. *A sub-sample $U_{m,n}^{(-i)}$ is valid with respect to a random walk with transition probability $\mathbb{P}[\cdot|\cdot]$ if $\mathbb{P}[U_{m,n}^{(i+1)}|U_{m,n}^{(i-1)}] > 0$ for $i \geq 2$ and $U_{m,n}^{(i+1)} = U_{m,n}^{(i)}$ for $i=1$, otherwise it is invalid.*

Definition 4 is general and it can be applied to any random walk algorithms. It states that a sub-sample $U_{m,n}^{(-i)}$ with respect to a random walk algorithm is valid if it can be a trajectory of this random walk with the same initial state $U_{m,n}^{(1)}$. For example, one valid sub-sample can be $U_{m,n}^{(-L_m)} = [U_{m,n}^{(1)}, \dots, U_{m,n}^{(L_m-1)}]$, which is obtained by deleting the last sample of $U_{m,n}$. Invalid sub-samples may cause the bias estimation (i.e., Equation (6)) being inaccurate. To relieve this problem, we design bootstrapping algorithms to estimate bias with sub-sample selection.

Bootstrapping via One Valid Sub-Sample. We first consider the sub-sample $U_{m,n}^{(-L_m)}$. This sub-sample has a nice statistical property that it is valid regardless of the samples in $U_{m,n}$. Namely, it is identically distributed as a random walk trajectory with length $L_m - 1$ and initial state $U_{m,n}^{(1)}$. This nice

statistical property of $U_{m,n}^{(-L_m)}$ is essential for analytical studies. Formally, we estimate the bias via the sub-sample $U_{m,n}^{(-L_m)}$ as

$$\widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n}) = (L_m - 1) \left(\hat{\Theta}_{m,n}^{(-L_m)} - \hat{\Theta}_{m,n} \right). \quad (7)$$

In Equation (7), we take the average over one sub-sample, i.e., $\hat{\Theta}_{m,n}^{(-L_m)}$, while we take the average of L_m sub-samples in Equation (6), i.e., $\sum_{i=1}^{L_m} \hat{\Theta}_{m,n}^{(-i)}/L_m$. Taking the average over one sub-sample leads to a larger variance of the bias estimation $\widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n})$ than that of $\widehat{\text{Bias}}(\hat{\Theta}_{m,n})$ with all the sub-samples (Equation (6)). Namely, our approach is at the cost of increasing the variance. As we have shown in Section 3, one can increase the number of parallel random walks to reduce the variance. Based on the new bias estimator $\widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n})$, we implement an oracle to bootstrap bias in Algorithm 5.

Algorithm 5. `BootBiasOracle`($U_{m,n}$, `RWOracle`) With One Valid Sub-Sample

- 1: $U_{m,n}^{(-L_m)} \leftarrow [U_{m,n}^{(1)}, \dots, U_{m,n}^{(L_m-1)}]$
- 2: $\hat{\Theta}_{m,n}^{(-L_m)} \leftarrow \text{RWOracle.Estimate}(U_{m,n}^{(-L_m)})$
- 3: $\hat{\Theta}_{m,n} \leftarrow \text{RWOracle.Estimate}(U_{m,n})$
- 4: Estimate the bias

$$\widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n}) \leftarrow (L_m - 1) \left(\hat{\Theta}_{m,n}^{(-L_m)} - \hat{\Theta}_{m,n} \right)$$

- 5: **return** $\widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n})$

To illustrate, consider $L_m = 3$. Then, we have $U_{m,n} = \{U_{m,n}^{(1)}, U_{m,n}^{(2)}, U_{m,n}^{(3)}\}$. We only have one sub-sample in consideration, i.e., $U_{m,n}^{(-3)} = [U_{m,n}^{(1)}, U_{m,n}^{(2)}]$. We can calculate the estimator for the sub-sample as: $\hat{\Theta}_{m,n}^{(-3)} = \text{RWOracle.Estimate}([U_{m,n}^{(1)}, U_{m,n}^{(2)}])$. Similarly, $\hat{\Theta}_{m,n} \leftarrow \text{RWOracle.Estimate}([U_{m,n}^{(1)}, U_{m,n}^{(2)}, U_{m,n}^{(3)}])$. The bias can be calculated as: $\widehat{\text{Bias}}(\hat{\Theta}_{m,n}) = 2(\hat{\Theta}_{m,n}^{(-3)} - \hat{\Theta}_{m,n})$. We next state theoretical guarantees for Algorithm 5.

Theorem 3. *Suppose the graph \mathcal{G} is connected and the T is bounded and expandable with $a_1(\mu) > 0$. Denote a space of all probability distribution over the attribute set \mathcal{X} as $\Omega = \{\rho | \sum_{y \in \mathcal{X}} \rho(y) = 1, \rho : \mathbb{V} \rightarrow [0, 1]\}$. Let $\Phi_{m,n}$ denote a probability distribution over Ω induced by $U_{m,n}$*

$$\Phi_{m,n}(\rho) = \mathbb{P} \left[\frac{1}{L_m} \sum_{i=1}^{L_m} \mathbb{1}_{\{x(U_{m,n}^{(i)})=y\}} = \rho(y) \right].$$

Let Ψ_n denote a probability distribution over Ω such that $\Psi_n(\rho) = \mathbb{P}[\mu_n = \rho]$. If $\|\Psi_{L_m} - \Phi_{m,n}\|_{TV} \leq o(1/L_m)$, where $\|\cdot\|_{TV}$ denotes the total variation distance. Algorithm 5 corrects the bias from $\text{Bias}(\hat{\Theta}_{m,n}) = O(1/L_m)$ to $\text{Bias}(\hat{\Theta}_{m,n}^{\text{vs}}) = O(1/L_m^2)$, where $\hat{\Theta}_{m,n}^{\text{vs}} \triangleq \hat{\Theta}_{m,n} - \widehat{\text{Bias}}_{\text{vs}}(\hat{\Theta}_{m,n})$ denotes the corrected estimator.

Theorem 3 states sufficient conditions under which the Algorithm 5 provides theoretical guarantees on the bias reduction. Note that in Theorem 3, the graph can be of any general topology provided that it is connected. It only requires an extra condition on the sample $U_{m,n}$ induced distribution over \mathcal{X} , which can be achieved by an appropriate

TABLE 1
Overall Statistics of Four Datasets

	# of nodes	# of edges	# of communities
com-Amazon	334,863	925,872	75,149
wiki-topcats	1,791,489	28,511,807	17,364
com-Orkut	3,072,441	117,185,083	6,288,363
com-LiveJournal	3,997,962	34,681,189	287,512

burn-in period. The following theorem states the computational complexity of Algorithm 5.

Theorem 4. *Under the $RWOracle$ implemented in Algorithms 2 or 3. The computational complexity of the bootstrapping oracle $BootBiasOracle(U_{m,n}, RWOracle)$ implemented in Algorithm 5 is $O(L_m)$.*

Compared with Theorem 2, Algorithm 5 reduces the computational complexity from quadratic in L_m to linear in L_m .

Bootstrapping via All Valid Sub-Samples. To reduce the variance of Algorithm 5, we propose an algorithm to estimate bias with all valid sub-samples. Due to page limit, we present them in our supplementary file, available online.

6 EXPERIMENT I: THE BIAS

We conduct experiments on four real-world datasets published on SNAP¹ to evaluate our algorithmic framework on bias reduction. Experiment results further confirm the superior performance of our algorithm in reducing the bias. *Due to page limit, some experiment results are presented in our supplementary file, available online.*

6.1 Experiment Setting

Dataset. Table 1 summarizes four datasets published on SNAP. We have two reasons in selecting them: (1) each node is associated with an attribute, i.e., community; (2) they are from four diverse applications. For each dataset, if it is a directed graph, we add reciprocal edges to make it an undirected one. In each dataset we consider two types of attributes: (1) the community associated with a node; (2) the degree of a node.

The scale of the graphs (i.e., millions of nodes) used in our experiment is widely used many previous works for the evaluations of random walk sampling algorithms [15], [32], [33], [34]. Note that our objective is using these graphs to compare our algorithm with other baseline algorithms. In academia, we may not have the chance to run experiments in the real-world setting, where calculating over the entire graph is impossible. The logic is that if we can show our algorithm has superior performance on these graphs, one may expect that our algorithm can have superior in the real-world setting.

Statistical Estimation Model. Consider the degree as the attribute of nodes. We have $\mathcal{X} = \{1, \dots, d_{\max}\}$ and $x(v)$ denotes the degree of node v . Furthermore, $\mu(y)$ denotes the fraction of nodes with degree $y \in \mathcal{X}$. We aim to estimate the standard deviation of degree

$$T(\mu) = \sqrt{\sum_{y \in \mathcal{X}} \mu(y)(y - \bar{y})^2}, \quad (8)$$

where $\bar{y} = \sum_{y \in \mathcal{X}} \mu(y)y$ denotes the average degree.

Consider the community as the attribute of nodes. We first rank the community ID, then based on the ranked list we divide the community ID into $K \in \mathbb{N}_+$ groups such that each group contains the same number of community ID. A node has attribute $k = 1, \dots, K$, if it belongs to a community with ID in group k . Thus, we have $\mathcal{X} = \{1, \dots, K\}$. We further set $\mu(k)$ as the fraction of nodes with community ID belonging to group k . Note that $\sum_{k=1}^K \mu(k) \neq 1$, as a node may belong to multiple communities or a node may not belong to any communities. Note that our framework applies to the case that $\sum_{k=1}^K \mu(k) \neq 1$. We estimate the variation of $\mu(k)$ as

$$T(\mu) = \left(\sum_{k=1}^K \frac{1}{K} \left| \mu(k) - \frac{\mu(1) + \dots + \mu(K)}{K} \right|^c \right)^{1/c}, \quad (9)$$

where $c \in \mathbb{R}_+$. When $c > 2$, the above statistic corresponds to generalized standard deviation.

Baseline & Parameter setting. To demonstrate the versatility of our framework, we apply it to reduce the bias of Metropolis random walk and simple random walk. When Metropolis random walk [29] serves as the baseline, we compare: (1) MR , which is a variant of Algorithm 2 without bias reduction, i.e., it is the Metropolis random walk; (2) JKM , which is a variant of Algorithm 2 and it uses Jackknife, i.e., Algorithm 4, to reduce the bias of Metropolis random walk; (3) VSM , which uses our bootstrapping oracle with valid sub-sample selection, i.e., Algorithm 5, to reduce bias of Metropolis random walk. When simple random walk serves as the baseline, we compare: (1) SRW , which extends Algorithm 2 to simple random walk, i.e., it is the simple random walk algorithm; (2) JKS , which extends Algorithm 4 to simple random walk, i.e., it uses Jackknife to reduce the bias of simple random walk; (3) VSS , which extends Algorithm 5 to simple random walk, i.e., it uses our bootstrapping oracle with valid sub-sample selection to reduce the bias of simple random walk.

We rank nodes based on the ID and then select ranked $(\lfloor \frac{1}{M} \rfloor \text{th}, \dots, (\lfloor \frac{M}{M} \rfloor \text{th})$ nodes as M initial points. Random walks in each group have the same length, i.e., $L_1 = \dots = L_M = L$. We also run the same number of parallel random walks on each initial point, i.e., $N_1 = \dots = N_M = N$. By default, we set the burn-in period to $\tilde{L} = 0$. We will vary \tilde{L} to study the impact of burn-in period on the effectiveness and efficiency of our method. Unless we state explicitly, we consider the following default parameters, i.e., $M=10, L=50, N=10^5, c=2$ and $K=5$ to compute the bias. Note that running this large number of random walks on each initial point is to ensure an accurate estimation of the bias via the Monte Carlo method.

6.2 Impact of Initial Points

We vary the number of initial points M from 4 to 18. We fix the total number of random walks to be $MN = 10^6$. Consider the case that community servers as node attribute. Fig. 1 shows the bias of estimating the statistic derived in Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the bias of Metropolis random walk MR by as high as 40%. This reduction ratio varies slightly as

1. <http://snap.stanford.edu/data/index.html>

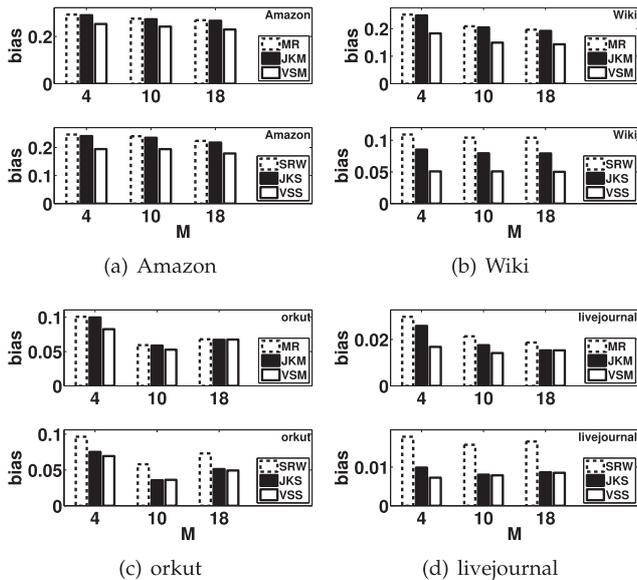


Fig. 1. Impact of M on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [*community as attribute*].

the number of initial points varies from $M = 4$ to $M = 18$. Furthermore, our bootstrapping method *VSM* reduces the bias of Jackknife method *JKM* by as high as 30%. Furthermore, our bootstrapping method *VSS* can reduce the bias of simple random walk *SRW* by as high as 50%. This reduction ratio varies slightly when the number of initial points varies from $M = 4$ to $M = 18$. When simple random walk serves as the baseline, our bootstrapping method *VSS* reduces the bias of Jackknife method *JKS* by as high as 30%.

Consider the case that node degree serves as node attribute. Fig. 2 shows the bias of estimating the statistic derived in Equation (8). One can observe that our bootstrapping method still reduces the bias of both Metropolis random walk and simple random walk significantly (as high as 80%) when the number of initial points varies from $M = 4$ to $M = 18$. Furthermore, our bootstrapping method reduces more bias than the Jackknife method in most cases when the number of initial points varies from $M = 4$ to $M = 18$.

Lessons Learned. Under different number of initial points, our bootstrapping method reduces the bias of both Metropolis random walk and simple random walk significantly (i.e., as high as 80%). It can also reduce the bias of the Jackknife method significantly in most cases, i.e., as high as 60%.

6.3 Impact of Sample Length

We vary the sample length L from 10 to 60, while set the other parameters as their default values. Consider the case that community serves as node attribute. Fig. 3 shows the bias of six algorithms (i.e., *MR*, *JKM*, *VSM*, *SRW*, *JKS* and *VSS*) evaluated on four datasets in Table 1, where the statistic under estimation is derived in Equation (9). One can observe that the bias of all these six algorithms decrease when the sample length L varies from 10 to 60. Namely, the bias decrease in sample length. Among *MR*, *JKM* and *VSM*, our *VSM* has the smallest bias followed by the *JKM*. Using a sample length of $L = 50$, our *VSM* reduces the bias of *MR* (or *JKM*) by as high as 40% (30%). Our bootstrapping method *VSS* reduces the bias of *SRW* (or *JKS*) by as high as 50% (30%).

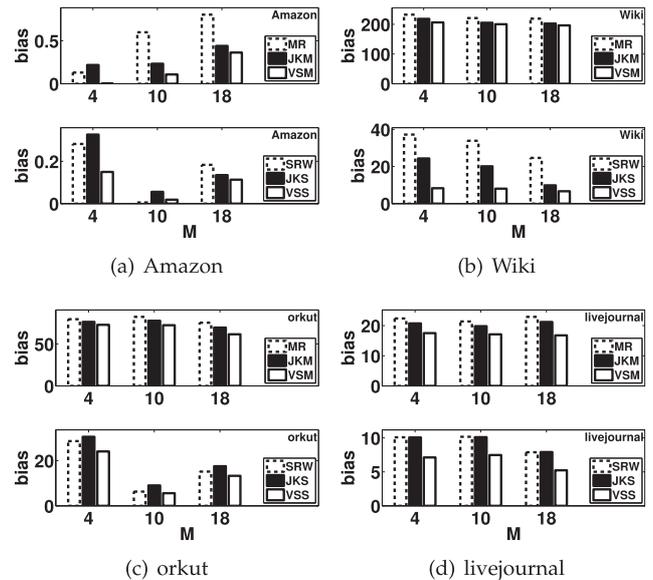


Fig. 2. Impact of M on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [*degree as attribute*].

Consider the case that node degree serves as node attribute. Fig. 4 shows the bias of estimating the statistic stated in Equation (8). One can observe that among *MR*, *JKM* and *VSM*, our *VSM* has the smallest bias. Using a sample length of $L=50$, *VSM* reduces the bias of *MR* (or *JKM*) by as high as 80% (50%). Furthermore, our bootstrapping method *VSS* reduces the bias of *SRW* (or *JKS*) by as high as 70% (60%).

Lessons Learned. Under different sample lengths, our bootstrapping method reduces the bias of both Metropolis random walk and simple random walk significantly and it can also reduce the bias of the Jackknife method.

6.4 Impact of Number of Community Groups

We vary the number of community groups K from 5 to 10. Fig. 5 shows the bias of estimating the statistic derived in

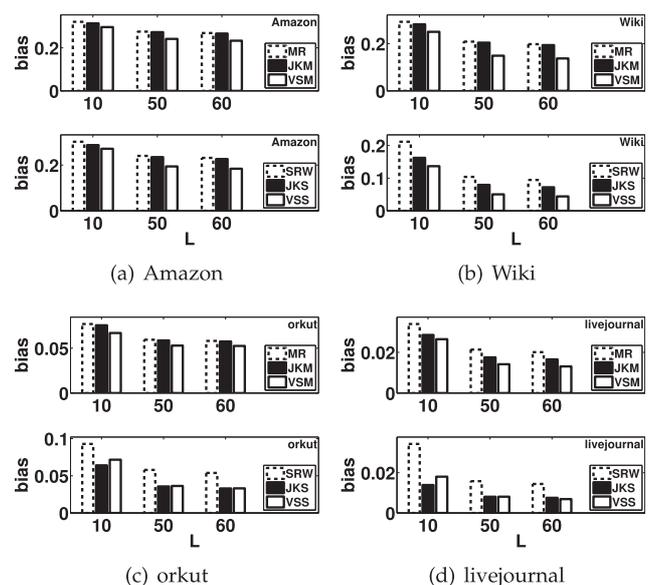


Fig. 3. Impact of sample length L on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [*community as attribute*].

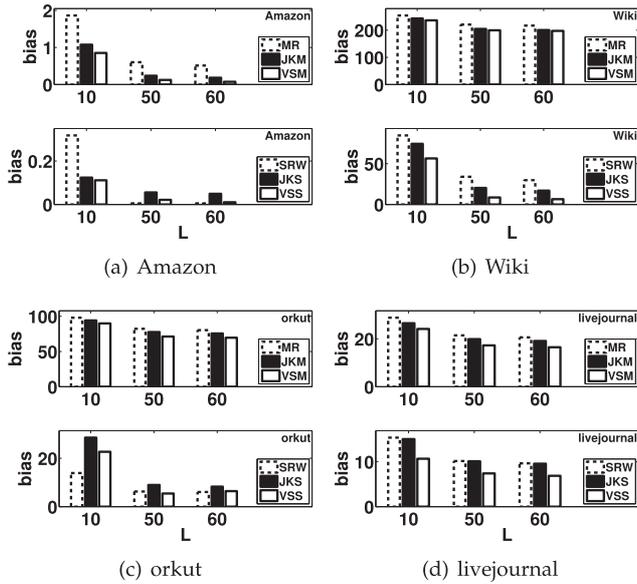


Fig. 4. Impact of sample length L on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [degree as attribute].

Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the bias of Metropolis random walk MR by as high as 30%. This reduction ratio varies slightly as the number of community groups varies from $K = 5$ to $K = 10$. Our bootstrapping method VSM reduces the bias of Jackknife method JKM by as high as 25%. Furthermore, our bootstrapping method VSS can reduce the bias of simple random walk SRW by as high as 50% and this reduction ratio varies slightly when the number of community groups varies from $K = 5$ to $K = 10$. When simple random walk serves as the baseline, our bootstrapping method VSS can reduce the bias of the Jackknife method JKS by as high as 30%.

Lessons Learned. Under different number of community groups, our bootstrapping method reduces the bias of both Metropolis random walk and simple random walk

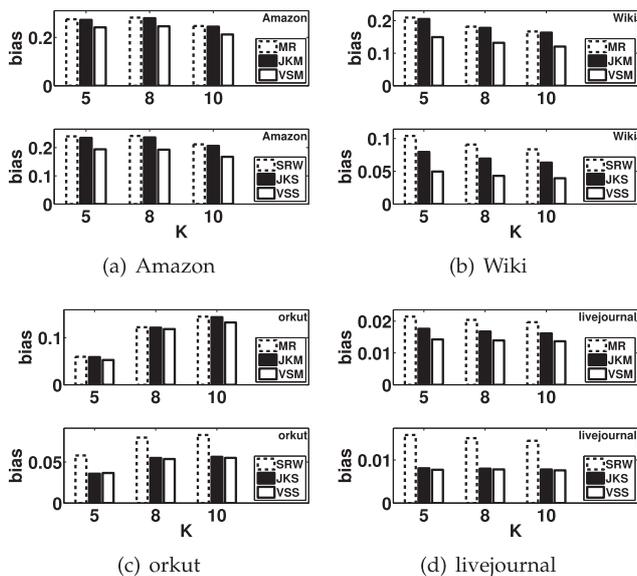


Fig. 5. Impact of K on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

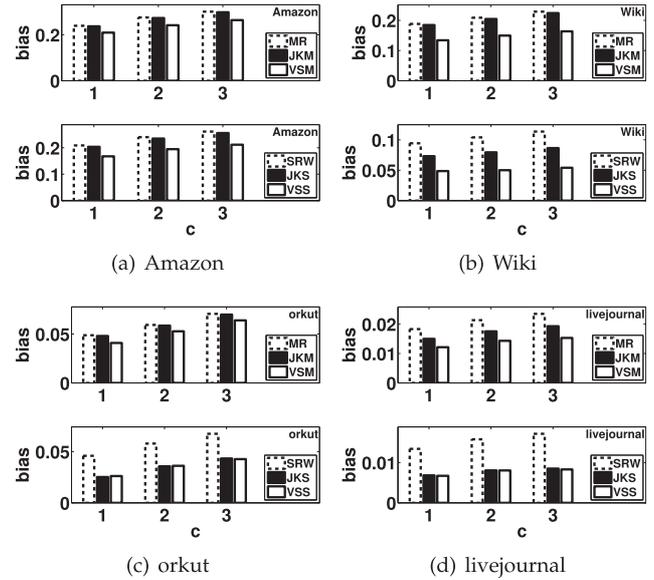


Fig. 6. Impact of statistical estimation model c on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

significantly (i.e., as high as 50%). It can also reduce more bias than the Jackknife method (i.e., as high as 30%).

6.5 Impact of Statistical Estimation Model

To study the impact of statistical estimation model on bias reduction, we vary the parameter c in the statistical estimation model derived in Equation (9) from 1 to 3. Fig. 6 shows the bias of estimating the statistic derived in Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the bias of Metropolis random walk MR by as high as 30%. This reduction ratio varies slightly as the parameter of the statistical estimation model varies from $c = 1$ to $c = 3$. When Metropolis random walk serves as the baseline, our bootstrapping method VSM reduces the bias of Jackknife method JKM by as high as 30%. Furthermore, our bootstrapping method VSS can reduce the bias of simple random walk SRW by as high as 50% and this reduction ratio varies slightly when c varies from 1 to 3. When simple random walk serves as the baseline, our bootstrapping method VSS can reduce the bias of the Jackknife method JKS by as high as 40%.

Lessons Learned. Under different statistical estimation models, our bootstrapping method reduces the bias of both Metropolis random walk and simple random walk significantly (by as high as 50%). It can also reduce the bias of Jackknife method by as high as 40%.

6.6 Impact of Burn-in Period

To study the impact of burn-in period on bias reduction, we vary the length of burn-in period \tilde{L} from 0 to 10000. Fig. 7 shows the bias of estimating the statistic derived in Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the bias of Metropolis random walk MR by as high as 25%. This reduction ratio varies slightly as the length of burn-in period \tilde{L} varies from 0 to 1000. When Metropolis random walk serves as the baseline, our bootstrapping method VSM reduces the bias of Jackknife method JKM by as high as 20%. Furthermore, our

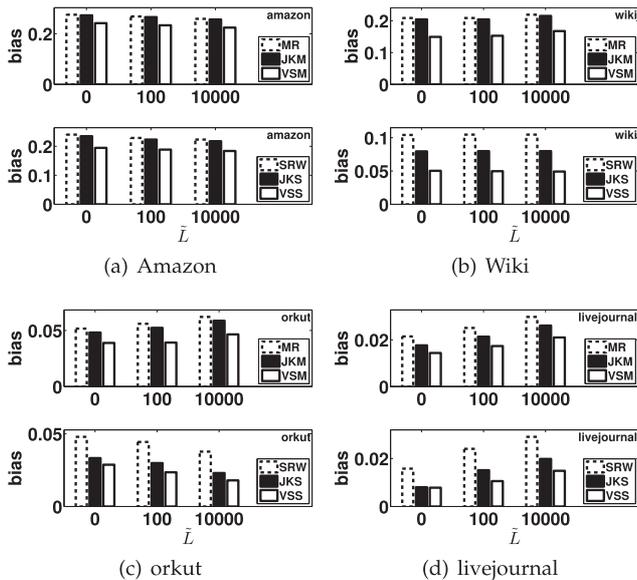


Fig. 7. Impact of length of burn-in period \tilde{L} on bias with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

bootstrapping method *VSS* can reduce the bias of simple random walk *SRW* by as high as 50% and this reduction ratio varies slightly when \tilde{L} varies from 0 to 10000. When simple random walk serves as the baseline, our bootstrapping method *VSS* can reduce the bias of the Jackknife method *JKS* by as high as 40%.

Lessons Learned. Under different length of burn-in period, our bootstrapping method reduces the bias of both Metropolis random walk and simple random walk significantly (by as high as 50%). It can also reduce the bias of Jackknife method by as high as 40%.

7 EXPERIMENT II: THE MEAN SQUARE ERROR

We conduct experiments to show how to use our bootstrapping algorithm to attain different trade-offs between the sample complexity and mean square error of an estimator. Experiment results further show that our bootstrapping method can reduce the mean square error of an estimator significantly by 1000 random walks. *Due to page limit, some experiment results are presented in supplementary file, available online.*

7.1 Experiment Setting

We consider the same experiment setting as Section 6, except that we study the mean square error of each algorithm when the total number of parallel random walks is small. We use Monte method to estimate the mean square error of each algorithm. In particular, we repeat each algorithm for 1000 times, and use the average of the outputs in these 1000 times to estimate the mean square error. Following previous works [9], [10], [33], [34], we consider the relative mean square error (RMSE), i.e., $\text{RMSE}(\hat{\Theta}) = \text{MSE}(\hat{\Theta})/\theta^2$, to eliminate the scale bias.

Due to page limit, we only consider the case that community serves as attribute. For the case of node degree serving as attribute, one can expect similar results, because the trade-off between the sample complexity and mean square

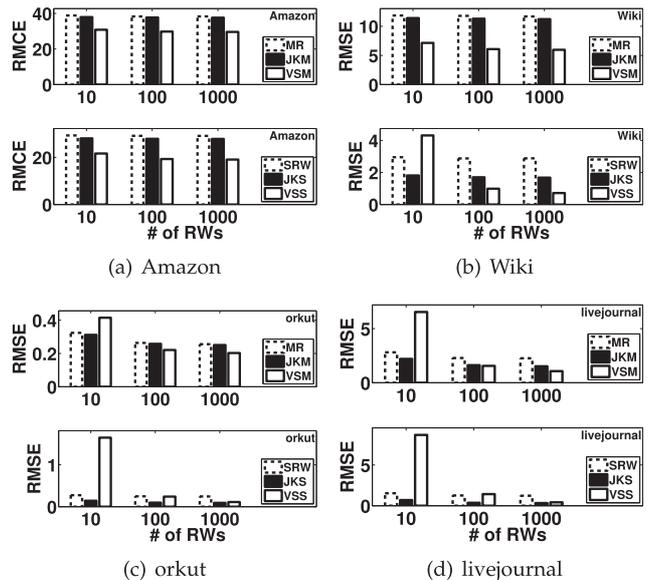


Fig. 8. Impact of number of random walks on the RMSE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

error is governed by the bias reduction and these two cases have similar bias reduction.

7.2 Impact of Number of Random Walkers

Note that the number of random walks equals MN . We set all the parameters except N as their default values stated in Section 6. We vary N from 1 to 100 such that the number of random walks MN varies from 10 to 1000. Fig. 8 shows the RMSE of six algorithms described in Section 6.1 evaluated on four datasets in Table 1, where the statistic in estimation is derived in Equation (9). One can observe that when the number of random walks is small, i.e., $MN = 10$, our *VSM* may have a larger RMSE than *MR* and *JKM* in some cases. This is because our *VSM* has a larger variance than *MR* and *JKM* and when the total number of random walks is small, the variance dominates. When the number of random walks is large, i.e., $MN = 1000$, our *VSM* reduces the RMSE of *MR* (*JKM*) by as high as 50% (40%). This is because when the total number of random walks is large, the bias dominates. Furthermore, when the number of random walks is small, i.e., $MN = 10$, our *VSS* may have a larger RMSE than *SRW* and *JKS* in some cases. When the number of random walks is large, i.e., $MN = 1000$, our *VSS* reduces the RMSE of *SRW* (*JKS*) by as high as 70% (50%).

Lessons Learned. Our bootstrapping method reduces the RMSE of Metropolis random walk and simple random walk significantly (i.e., as high as 70%) by no more than one thousand random walks and it can also reduce the bias of Jackknife method by as high as 50%.

7.3 Impact of Sample Length

We vary the sample length L from 10 to 60. We set the number of random walks to be 1000. All the other parameters are set as default values. Fig. 9 shows the RMSE of six algorithms described in Section 6.1 under the statistic derived in Equation (9). One can observe that the RMSE of these six algorithms decreases as the sample length L increases.

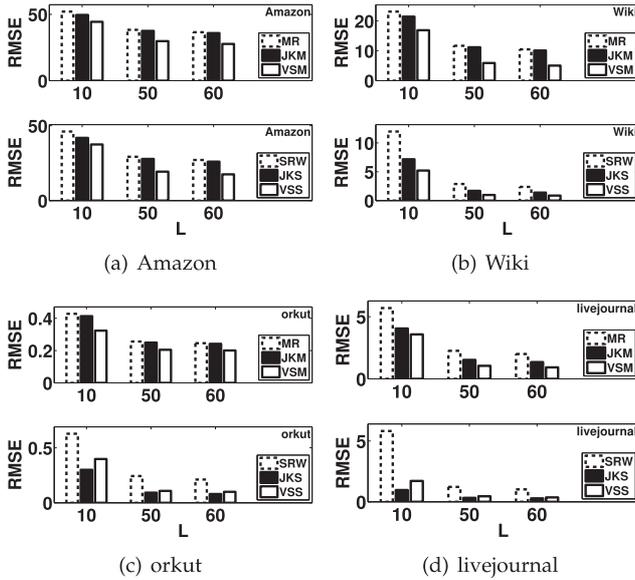


Fig. 9. Impact of walk length L on the RMSE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

Furthermore, our VSM reduces the RMSE of MR and JKM by as high as 50%. Lastly, our VSS reduces the RMSE of SRW (JKS) by as high as 60% (40%).

7.4 Impact of Number of Community Groups

We consider the same setting as Section 7.3, except we set the sample length L to be 50 and vary the number of community groups K from 5 to 10. Fig. 10 shows the RMSE of six algorithms under the statistic derived in Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the RMSE of Metropolis random walk MR by as high as 50%. This reduction ratio varies slightly as we vary the number of community groups from $K = 5$ to $K = 10$. Our VSM reduces the RMSE of Jackknife method JKM by as high

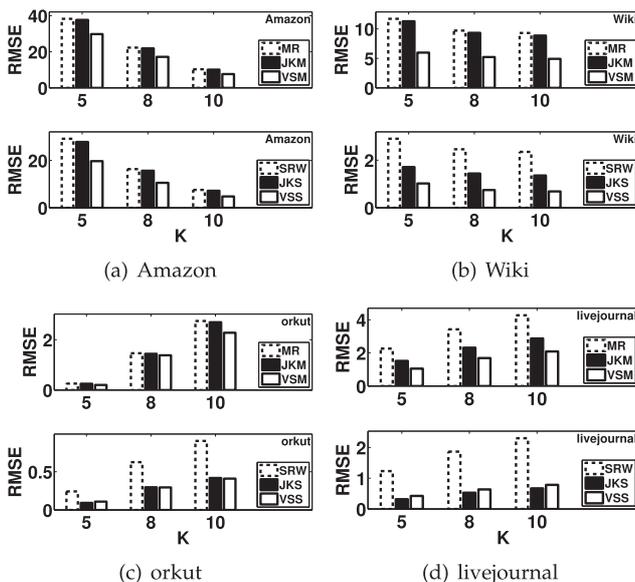


Fig. 10. Impact of community groups K on the RMSE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

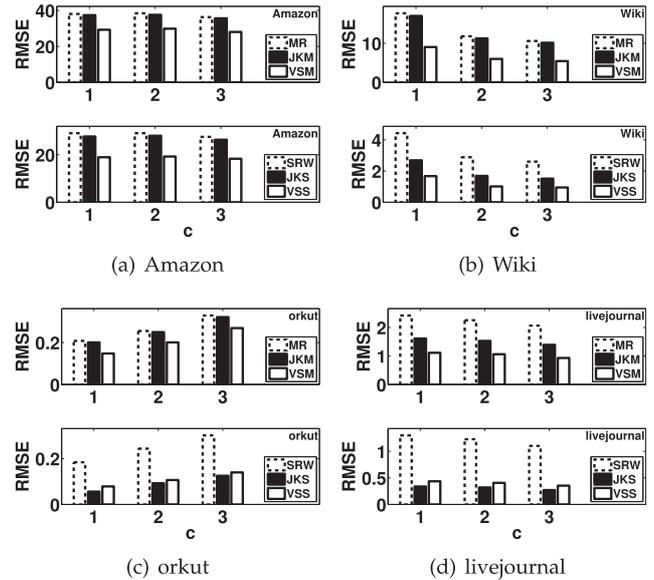


Fig. 11. Impact of statistical estimation model c on the RMSE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

as 40%. Furthermore, our bootstrapping method VSS can reduce the RMSE of simple random walk SRW by as high as 70% and this reduction ratio varies slightly when the number of community groups varies from $K = 5$ to $K = 10$. Lastly, when simple random walk serves as the baseline, our VSS can reduce the RMSE of the Jackknife method JKS by as high as 50%.

Lessons Learned. Under different number of community groups and one thousand random walks, our bootstrapping method reduces the RMSE of both Metropolis random walk and simple random walk significantly (i.e., as high as 70%) and it can also reduce the RMSE of Jackknife method by as high as 50%.

7.5 Impact of Statistical Estimation Model

We vary the parameter c of Equation (9) from 1 to 3. All the other parameters are set the same as Section 7.3, except the sample length L is set to be 50. Fig. 11 shows the RMSE of six algorithms under the statistic derived in Equation (9). One can observe that our bootstrapping algorithm VSM can reduce the RMSE of Metropolis random walk MR by as high as 50%. This reduction ratio varies slightly as we vary the parameter of the statistical estimation model from $c = 1$ to $c = 3$. Our bootstrapping method VSM reduces the RMSE of Jackknife method JKM by as high as 50%. Furthermore, our bootstrapping method VSS can reduce the RMSE of simple random walk SRW by as high as 70% and this reduction ratio varies slightly when c varies from 1 to 3. Lastly, when simple random walk serves as the baseline, our bootstrapping method VSS reduces the RMSE of the Jackknife method JKS by as high as 40%.

Lessons Learned. Under different statistical estimation models and one thousand random walks, our bootstrapping method reduces the RMSE of both Metropolis random walk and simple random walk significantly (by as high as 70%) and it can also reduce the RMSE of Jackknife method by as high as 40%.

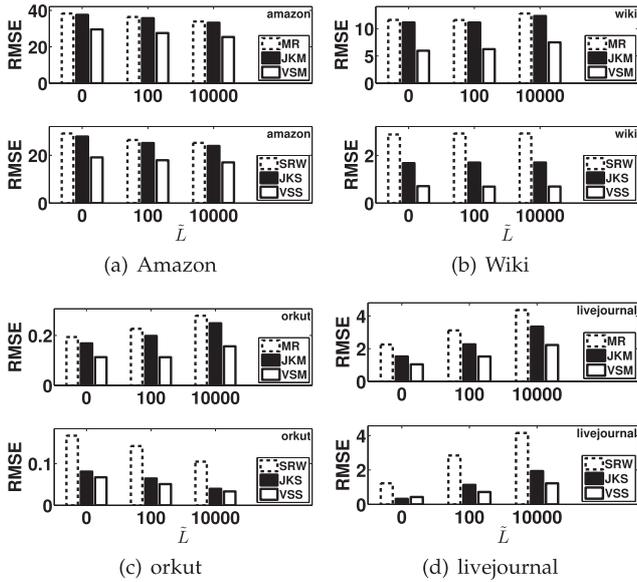


Fig. 12. Impact of length of burn-in \tilde{L} on RMSE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

7.6 Impact of Burn-in Period

We vary the length of burn-in period \tilde{L} from 0 to 10000. All the other parameters are set the same as Section 7.3, except the sample length L is set to be 50. Fig. 12 shows the RMSE of six algorithms under the statistic derived in Equation (9). One can observe that our bootstrapping algorithm *VSM* can reduce the RMSE of Metropolis random walk *MR* by as high as 50%. This reduction ratio varies slightly as we vary \tilde{L} from 0 to 10000. Our bootstrapping method *VSM* reduces the RMSE of Jackknife method *JKM* by as high as 50% as well. Furthermore, our bootstrapping method *VSS* can reduce the RMSE of simple random walk *SRW* by as high as 80% and this reduction ratio varies slightly when \tilde{L} from 0 to 10000. Lastly, when simple random walk serves as the baseline, our bootstrapping method *VSS* reduces the RMSE of the Jackknife method *JKS* by as high as 60%.

Lessons Learned. Under different length of burn-in period and one thousand random walks, our bootstrapping method reduces the RMSE of both Metropolis random walk and simple random walk significantly (by as high as 80%) and it can also reduce the RMSE of Jackknife method by as high as 60%.

8 EXPERIMENT III: BEYOND MSE

In this section, we evaluate our proposed algorithms beyond the mean square error metric. In particular, we evaluate our proposed algorithms under the relative mean absolute error (RMAE) metric, i.e., $\text{RMAE}(\hat{\Theta}) \triangleq \mathbb{E}[|\hat{\Theta} - \theta|/|\theta|]$, and the relative mean cubic absolute error (RCAE), i.e., $\text{RCAE}(\hat{\Theta}) \triangleq \mathbb{E}[|\hat{\Theta} - \theta|^3]/|\theta|^3$. We consider the same experiment setting as Section 7. Due to page limit, some experiment results are presented in our supplementary file, available online.

8.1 Impact of Number of Random Walks

Note that the number of random walks equals MN . We set all the parameters except N as their default values stated in

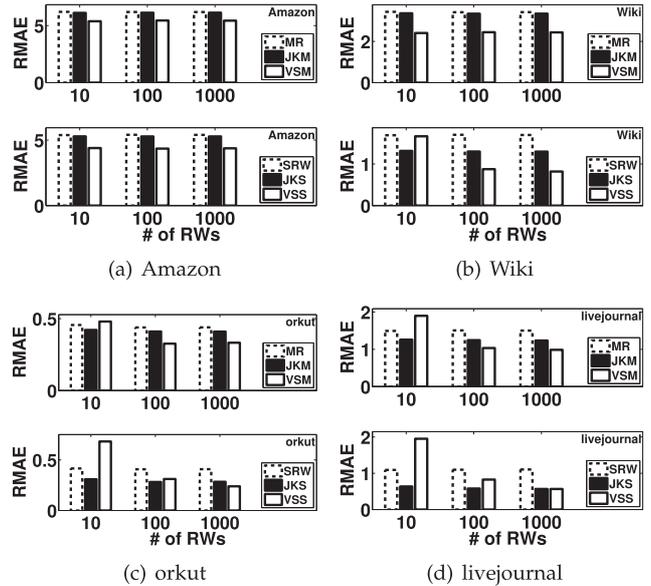


Fig. 13. Impact of number of random walks on the RMAE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

Section 6. We vary N from 1 to 100 such that the number of random walks MN varies from 10 to 1000. Fig. 13 shows the RMAE of six algorithms evaluated on four datasets in Table 1, where the statistic in estimation is derived in Equation (9). One can observe that when the number of random walks is small, i.e., $MN = 10$, our *VSM* may have a larger RMAE than *MR* and *JKM* in some cases. This is because our *VSM* has a larger variance than *MR* and *JKM* and when the total number of random walks is small, the variance dominates. When the number of random walks is large, i.e., $MN = 1000$, our *VSM* reduces the RMAE of *MR* (*JKM*) by as high as 50% (40%). This is because when the total number of random walks is large, the bias dominates. Furthermore, when the number of random walks is small, i.e., $MN = 10$, our *VSS* may have a larger RMAE than *SRW* and *JKS* in some cases. When the number of random walks is large, i.e., $MN = 1000$, our *VSS* reduces the RMAE of *SRW* (*JKS*) by as high as 60% (30%). Similar statements holds for the RMCE metric, as shown in Fig. 14. In summary, our proposed algorithms still have superior performance over baselines under both the RMAE and RMCE metric.

9 RELATED WORK

Random Walk. Random walk sampling is a mainstream method to generate representation samples from large scale graphs [16]. Two fundamental random walk sampling algorithms are (1) the simple random walk [30], and (2) the Metropolis random walk [29]. A number of variants of random walk sampling algorithms were proposed to improve the estimation accuracy. Rasti *et al.* [32] proposed an algorithm which incorporate respondent-driven sampling into Metropolis random walk. Ribeiro *et al.* [33] developed a coordinated multidimensional random walk sampling algorithm. Kurant *et al.* [9] proposed a stratified weighted random walk algorithm. Jin *et al.* [21] and Xu *et al.* [34] proposed random walk sampling algorithms with jumps. Lu *et al.* [3] developed an algorithm to approximate the bias

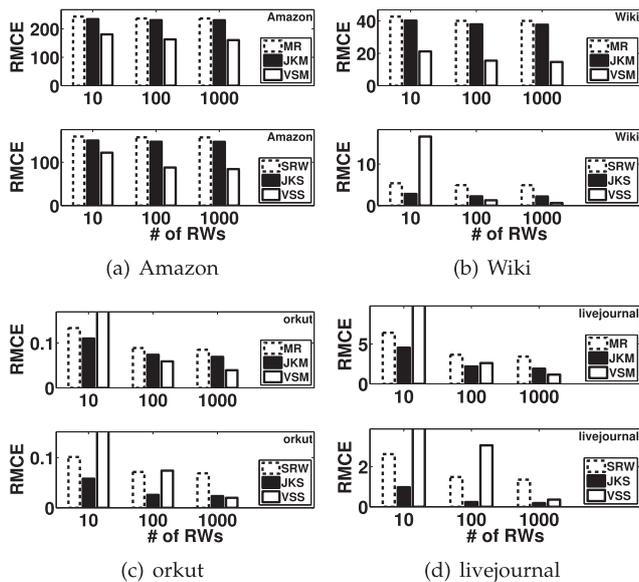


Fig. 14. Impact of number of random walks on the RMCE with Metropolis random walk (MR) and simple random walk (SRW) as baselines. [community as attribute].

of estimating the population size via random walk. Lee *et al.* [10] proposed two algorithms to reduce the asymptotic variance of estimators: (1) non-backtracking random walk and (2) random walk with delayed acceptance. Li *et al.* [20] further extended the delayed acceptance method to Metropolis random walk, which is shown to further reduce the asymptotic variance. Lu *et al.* [35] showed that the harmonic mean estimator for average degree can reduce the estimation variance significantly. Zhou *et al.* [8], [28], [36] proposed history dependent random walk sampling algorithms, which are shown to have a fast convergence speed. Li *et al.* [15] further improved them by considering the walking history and next-hop candidates. Essentially, all these algorithms improve estimation accuracy via designing random walk strategies, i.e., they focus on generating samples. These algorithms only use simple average to do the estimation. Our framework is orthogonal to their works and complements them. First, we focus on how to utilize samples to produce accurate estimation. In particular, we achieve this by applying bootstrapping techniques to design a bias reduction algorithm. Second, our framework utilizes recent graph processing systems which can run millions of random walks on a consumer-level personal computer. Third, our framework is generic and it can be applied to improve the estimation accuracy of any random walk strategies.

Our work is orthogonal to works on parallel random walk algorithms and systems [22], [23], [24], [25], [26] in the sense that it is independent of the design of parallel random walk algorithms, i.e., it can be applied to a broad class of parallel random walk algorithms.

Bootstrapping. Bootstrapping is a technique for statistical estimation. There are a variety of bootstrapping techniques such as Efron bootstrap and Jackknife for different applications or settings [27], [31], [37], [38], [39]. Model-based bootstrapping techniques require a large number of samples to accurately reconstruct the model. Thus it is not suitable for random walk sampling applications because the sample

size is usually small. Model-free bootstrapping techniques such as Jackknife have the issue of invalid sub-samples in random walk sampling applications. Our work proposes a new variant of the Jackknife method with sub-sample selection which is fine tuned for the random walk sampling algorithms. In our framework we allow the variance to increase whenever the bias can be further reduced. We overcome the increased variance by parallel random walks.

10 CONCLUSION

This paper develops an algorithmic framework to improve the accuracy of random walk based statistical estimation over graphs. We apply the bootstrapping technique to design a bias reduction algorithm. Our algorithmic framework enables one to attain different trade-offs between the sample complexity and the error of statistical estimation. Also, our bias reduction algorithm is generic and can be applied to optimize a large class of random walk sampling algorithms. We provide theoretical guarantees and computational complexity analysis of our proposed bias reduction algorithms. Extensive experiments on four public datasets confirm the effectiveness and efficiency of our proposed algorithmic framework under the mean square metric and beyond.

REFERENCES

- [1] M. Kim and J. Leskovec, "Modeling social networks with node attributes using the multiplicative attribute graph model," in *Proc. 27th Conf. Uncertainty Artif. Intell.*, 2011, pp. 400–409.
- [2] G. B. Giannakis, Y. Shen, and G. V. Karanikolas, "Topology identification and learning over graphs: Accounting for nonlinearities and dynamics," *Proc. IEEE*, vol. 106, no. 5, pp. 787–807, May 2018.
- [3] J. Lu and D. Li, "Bias correction in a small sample from big data," *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 11, pp. 2658–2663, Nov. 2013.
- [4] K. Nakajima and K. Shudo, "Estimating properties of social networks via random walk considering private nodes," in *Proc. 26th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2020, pp. 720–730.
- [5] X. Yang, H. Steck, Y. Guo, and Y. Liu, "On top-k recommendation using social networks," in *Proc. 6th ACM Conf. Recommender Syst.*, 2012, pp. 67–74.
- [6] Z. Xu, Y. Ke, Y. Wang, H. Cheng, and J. Cheng, "A model-based approach to attributed graph clustering," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2012, pp. 505–516.
- [7] S. Zhang, J. Yang, and V. Cheedella, "Monkey: Approximate graph mining based on spanning trees," in *Proc. IEEE 23rd Int. Conf. Data Eng.*, 2007, pp. 1247–1249.
- [8] Z. Zhou, N. Zhang, Z. Gong, and G. Das, "Faster random walks by rewiring online social networks on-the-fly," *ACM Trans. Database Syst.*, vol. 40, no. 4, pp. 1–36, 2016.
- [9] M. Kurant, M. Gjoka, C. T. Butts, and A. Markopoulou, "Walking on a graph with a magnifying glass: Stratified sampling via weighted random walks," in *Proc. ACM SIGMETRICS Joint Int. Conf. Meas. Model. Comput. Syst.*, 2011, pp. 281–292.
- [10] C.-H. Lee, X. Xu, and D. Y. Eun, "Beyond random walk and metropolis-hastings samplers: Why you should not backtrack for unbiased graph sampling," *ACM SIGMETRICS Perform. Eval. Rev.*, vol. 40, no. 1, pp. 319–330, 2012.
- [11] X. Xu, C.-H. Lee, and D. Y. Eun, "Challenging the limits: Sampling online social networks with cost constraints," in *Proc. IEEE Conf. Comput. Commun.*, 2017, pp. 1–9.
- [12] S. Agarwal, "Ranking on graph data," in *Proc. 23rd Int. Conf. Mach. Learn.*, 2006, pp. 25–32.
- [13] C. E. Priebe, D. L. Sussman, M. Tang, and J. T. Vogelstein, "Statistical inference on errorfully observed graphs," *J. Comput. Graph. Statist.*, vol. 24, no. 4, pp. 930–953, 2015.
- [14] D. Zhou and B. Schölkopf, "A regularization framework for learning from graph data," in *Proc. Workshop Statist. Relational Learn. Int. Conf. Mach. Learn.*, 2004, pp. 132–137.

- [15] Y. Li *et al.*, "Walking with perception: Efficient random walk sampling via common neighbor awareness," in *Proc. IEEE 35th Int. Conf. Data Eng.*, 2019, pp. 962–973.
- [16] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou, "Walking in Facebook: A case study of unbiased sampling of OSNs," in *Proc. IEEE INFOCOM*, 2010, pp. 1–9.
- [17] X. Chen, Y. Li, P. Wang, and J. C. Lui, "A general framework for estimating graphlet statistics via random walk," *Proc. VLDB Endowment*, vol. 10, pp. 253–264, 2016.
- [18] P. Wang *et al.*, "MOSS-5: A fast method of approximating counts of 5-node graphlets in large graphs," *IEEE Trans. Knowl. Data Eng.*, vol. 30, no. 1, pp. 73–86, Jan. 2018.
- [19] J. Zhao, P. Wang, J. C. Lui, D. Towsley, and X. Guan, "Sampling online social networks by random walk with indirect jumps," *Data Mining Knowl. Discov.*, vol. 33, no. 1, pp. 24–57, 2019.
- [20] R.-H. Li, J. X. Yu, L. Qin, R. Mao, and T. Jin, "On random walk based graph sampling," in *Proc. IEEE 31st Int. Conf. Data Eng.*, 2015, pp. 927–938.
- [21] L. Jin *et al.*, "Albatross sampling: Robust and effective hybrid vertex sampling for social graphs," in *Proc. 3rd ACM Int. Workshop MobiArch*, 2011, pp. 11–16.
- [22] A. Kyröla, "DrunkardMob: Billions of random walks on just a PC," in *Proc. 7th ACM Conf. Recommender Syst.*, 2013, pp. 257–264.
- [23] K. Vora, G. Xu, and R. Gupta, "Load the edges you need: A generic I/O optimization for disk-based graph processing," in *Proc. USENIX Conf. Usenix Annu. Tech. Conf.*, 2016, pp. 507–522.
- [24] H. Liu and H. H. Huang, "Graphene: Fine-grained IO management for graph computing," in *Proc. 15th Usenix Conf. File Storage Technol.*, 2017, pp. 285–299.
- [25] K. Yang, M. Zhang, K. Chen, X. Ma, Y. Bai, and Y. Jiang, "KnightKing: A fast distributed graph random walk engine," in *Proc. 27th ACM Symp. Oper. Syst. Princ.*, 2019, pp. 524–537.
- [26] R. Wang, Y. Li, H. Xie, Y. Xu, and J. C. Lui, "GraphWalker: An I/O-efficient and resource-friendly graph analytic system for fast and scalable random walks," in *Proc. USENIX Conf. Usenix Annu. Tech. Conf.*, 2020, Art. no. 38.
- [27] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*. Boca Raton, FL, USA: CRC Press, 1994.
- [28] Z. Zhou, "Faster sampling over online social networks," PhD dissertation, George Washington Univ., Washington, DC, USA, 2015.
- [29] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, "Equation of state calculations by fast computing machines," *J. Chem. Phys.*, vol. 21, no. 6, pp. 1087–1092, 1953.
- [30] L. Lovász *et al.*, "Random walks on graphs: A survey," *Combinatorics, Paul Erdos is Eighty*, vol. 2, no. 1, pp. 1–46, 1993.
- [31] B. Efron, *The Jackknife, the Bootstrap, and Other Resampling Plans*, vol. 38. Philadelphia, PA, USA: SIAM, 1982.
- [32] A. H. Rasti, M. Torkjazi, R. Rejaie, N. Duffield, W. Willinger, and D. Stutzbach, "Respondent-driven sampling for characterizing unstructured overlays," in *Proc. IEEE INFOCOM*, 2009, pp. 2701–2705.
- [33] B. Ribeiro and D. Towsley, "Estimating and sampling graphs with multidimensional random walks," in *Proc. 10th ACM SIGCOMM Conf. Internet Meas.*, 2010, pp. 390–403.
- [34] X. Xu, C.-H. Lee, and D. Y. Eun, "A general framework of hybrid graph sampling for complex network analysis," in *Proc. IEEE Conf. Comput. Commun.*, 2014, pp. 2795–2803.
- [35] J. Lu and H. Wang, "Variance reduction in large graph sampling," *Inf. Process. Manage.*, vol. 50, no. 3, pp. 476–491, 2014.
- [36] Z. Zhou, N. Zhang, and G. Das, "Leveraging history for faster sampling of online social networks," *Proc. VLDB Endowment*, vol. 8, pp. 1034–1045, 2015.
- [37] A. C. Davison and D. V. Hinkley, *Bootstrap Methods and Their Application*, vol. 1. Cambridge, U.K.: Cambridge Univ. Press, 1997.
- [38] J. Shao and D. Tu, *The Jackknife and Bootstrap*. Berlin, Germany: Springer, 2012.
- [39] S. N. Lahiri, *Resampling Methods for Dependent Data*. Berlin, Germany: Springer, 2013.



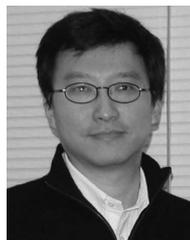
Hong Xie received the BEng degree from the School of Computer Science and Technology, University of Science and Technology of China, Hefei, China, and the PhD degree from the Department of Computer Science and Engineering, Chinese University of Hong Kong, Hong Kong, advised by Prof. John C. S. Lui. He is currently a research professor with the College of Computer Science, Chongqing University. He was a postdoctoral research fellow with CUHK and NUS.



Pei Yi received the BEng degree from Hohai University, Nanjing, China. He is currently working toward the master's degree in the College of Computer Science, Chongqing University, Chongqing, China, under the supervision of Prof. Hong Xie. His research interests include graph analytics, Markov Chain Monte Carlo, etc.



Yongkun Li received the bachelor's degree in computer science from the University of Science and Technology of China, Hefei, China, and the PhD degree in computer science from the Chinese University of Hong Kong, Hong Kong, advised by Prof. John C. S. Lui. He is an associate professor with the School of Computer Science and Technology, University of Science and Technology of China. Before that, he was a postdoctoral fellow with the Institute of Network Coding, working with Prof. John C. S. Lui and Prof. Patrick P. C. Lee.



John C. S. Lui (Fellow, IEEE) received the PhD degree in computer science from the University of California at Los Angeles, Los Angeles, California. He was a chairman with the CSE Department from 2005 to 2011. He is currently the Choh-Ming Li chair professor with the Department of Computer Science and Engineering, Chinese University of Hong Kong. His current research interests include communication networks, system security (e.g., cloud security, mobile security, etc.), network economics, network sciences, large-scale distributed systems, and performance evaluation theory. He is an elected member of the IFIP WG 7.3, and a Croucher senior research fellow. He was a recipient of the various departmental teaching awards and the CUHK Vice-Chancellors Exemplary Teaching Award.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.