

# On the Estimation Accuracy of Degree Distributions from Graph Sampling

Bruno Ribeiro, Don Towsley

**Abstract**—Estimating characteristics of large graphs via sampling is a vital part of the study of complex networks. In this work we present an in-depth study of the Mean Squared Error (MSE) of sampling methods such as independent random vertex (RV) and random edge (RE) sampling and crawling methods such as random walks (RWs), a.k.a. RDS, and the a Metropolis-Hastings algorithm whose target distribution is to uniformly sample vertices (MHRWu). This paper provides an upper bound for the MSE of a stationary RW as a function of the MSE of RE and the absolute value of the second most dominant eigenvalue of the RW transition probability matrix. We see that RW and RV sampling are optimal in respect to different weighted MSE optimizations and show when RW is preferable to RV sampling. Finally, we present an approximation to the MHRWu MSE. We evaluate the accuracy of our approximations and bound using simulations on large real world graphs.

## I. INTRODUCTION

A number of recent studies [2], [4], [7], [8], [12], [13], [15], [21], [16], [17], [22] (to cite a few) are dedicated to the characterization of network graphs. This paper represents a network as an undirected graph with labeled vertices and edges. Network characteristics of interest include the degree distribution, the average number of copies of a file in a peer-to-peer (P2P) network [8], [21], the assortativity coefficient [18], or the global clustering coefficient [18].

Characterizing graphs requires querying vertices and/or edges; each query has an associated resource cost (time, bandwidth, money). Querying the whole graph is often too costly. As a result, researchers have turned their attention to the estimation of graph characteristics based on incomplete (sampled) data.

*RV sampling:* In networks where each vertex is assigned a unique user-id (e.g., travelers and their passport numbers, Facebook, MySpace, Flickr, and Livejournal) a widespread practice is to perform random vertex (RV) sampling by querying randomly generated user-ids. However, uniform RV sampling may be undesirable when the user-id space is sparsely populated (in MySpace the ratio between the number valid users retrieved and the total number of queries is 10% [17]). Moreover, queries are often subject to

resource constraints (e.g., queries are rate-limited in Flickr, Livejournal [15], and Bittorrent [11]). As we see in this work, even when RV sampling is not severely resource-constrained, some characteristics may be better estimated with other sampling methods (e.g., the tail of the degree distribution of a graph).

*RE sampling:* In independent Random Edge (RE) sampling, a vertex is sampled by first sampling an edge independent and uniformly from the set of edges, and then randomly choosing one of the edge end points. In practice one should use both end points of a sampled edge. However, in order to simplify our analysis, we consider just one sampled vertex for each sampled edge. In real world networks, randomly sampling edges can be harder than randomly sampling vertices. Edges are not often associated to unique IDs that can be queried and online social networks such as Facebook, Twitter, MySpace, Livejournal, and Flickr, among others, do not provide an API that allows randomly sampling of edges.

*RW sampling:* An alternative, and often cheaper, way to sample a network is by means of a random walk (RW). RW sampling is preferred to other types of graph crawling, such as the breadth-first crawling used in [15], as one can obtain asymptotically unbiased estimates of a number of graph characteristics such as fraction of vertices with a given label [22], the degree distribution [22], and, more recently, assortativity and global clustering coefficients [18]. A RW samples a graph by moving a particle (walker) from a vertex to a neighboring vertex (over an edge). The probability by which the walker selects the next neighboring vertex determines the probability by which vertices and edges are sampled. We denote *standard RW* or just *RW* a random walk that sample neighbors *uniformly*. A Metropolis-Hastings walker, as seen later, selects the next neighboring vertex using a different rule. RWs are popular for sampling networks [7], [16], [22] in order to estimate their characteristics. One of the reasons behind the popularity of RW sampling is that it does not query invalid users unlike RV sampling.

*MHRW sampling:* The Metropolis-Hastings Random Walk (MHRW) is an accept-reject random walk-based sampling process that samples vertices according to a target distribution  $\gamma$ . In this work we are mostly interested in a

The authors are with the Computer Science Department, University of Massachusetts Amherst, Amherst, MA, 01003 USA e-mail: {ribeiro,towsley}@cs.umass.edu

MHRW that samples vertices uniformly, which we denote MHRWu. MHRWu have been used to uniformly sample peers in peer-to-peer networks [21] and Web pages [9]. Unfortunately, MHRWu is empirically known to have large estimation errors compared to RW estimates [7].

### Contributions

This paper presents the following contributions:

- 1) In Section III we prove that the Mean Squared Error (MSE) obtained by a stationary sequence of  $n$  RW sampled vertices is upper bounded by the MSE of  $n$  RE sampled vertices divided by  $(1 - \alpha)$ , where  $\alpha$  is the absolute value of the second most dominant eigenvalue of the RW transition probability matrix.
- 2) We present the graph sampling problem as the minimization of a weighted MSE sum. We illustrate our approach using the degree distribution as an example. For the degree distribution we see that RW sampling minimizes the MSE whose weights are the vertex degree squared (i.e., the weights give more importance to large degree vertices); while RV sampling minimizes the MSE with equal weights.
- 3) RW estimates have been observed to be more accurate than estimates obtained by MHRWu [7], [16]. We study how the Metropolis-Hastings mechanism tends to induce larger estimation errors than RW and RV sampling.

### Outline

The outline of this work is as follows. Section II presents definitions used in this paper. Section III presents an upper bound of the MSE of RW sampling as a function of the MSE of RE sampling and  $\alpha$ , the absolute value of the second most dominant eigenvalue of the RW transition probability matrix. In Section IV we present the graph sampling problem as minimizing the weighted MSE sum. In Section V we study how the Metropolis-Hastings mechanism tends to induce larger estimation errors than RW or even RV sampling. Section VII present simulation results that help corroborate our theoretical analysis. And finally Section IX presents our conclusions.

## II. DEFINITIONS

Let  $G = (V, E)$  be an undirected connected non-bipartite graph and let  $d_a$ ,  $a \in V$ , be the degree of vertex  $a$ . We denote  $\text{vol}(V) \triangleq \sum_{v \in V} d_v$ . We want to estimate

$$F = \sum_{v \in V} f(v). \quad (1)$$

from a sequence of vertices sampled from  $G$ . Let  $(Z_1, \dots, Z_n)$  be a stationary sequence of  $n$  sampled vertices, where  $P[Z_t = v] = \beta_v > 0$ ,  $\forall v \in V$ ,  $t = 1, \dots, n$ . Then

$$\hat{F}(Z_1, \dots, Z_n) \triangleq \frac{1}{n} \sum_{t=1}^n \frac{f(Z_t)}{\beta_{Z_t}}, \quad Z_i \in V, i = 1, \dots, n. \quad (2)$$

is an unbiased estimate of eq.(1). The estimator  $\hat{F}$  in eq.(2) is widely used to estimate  $F$ , see [18], [22] and the references therein.

The Mean Squared Error (MSE) of  $\hat{F}(Z_1, \dots, Z_n)$  is

$$E[(\hat{F}(Z_1, \dots, Z_n) - F)^2] = \text{var}(\hat{F}(Z_1, \dots, Z_n)), \quad (3)$$

as  $E[\hat{F}(Z_1, \dots, Z_n)] = F$ .

### III. A TIGHT UPPER BOUND OF THE RW ESTIMATION ERROR

Random walk (RW) sampling and random edge (RE) sampling are closely related. Let  $\pi = (\pi_a : a \in V)$ , denote the steady state probability distribution of the RW. A RW is time reversible, i.e.,  $\pi_v/d_v = \pi_u/d_u$  [14]. A consequence of time reversibility is that edges are sampled with equal probability,  $1/|E|$ . Thus, RW and RE differ only that edges sampled by a RW samples are correlated.

In what follows we present a tight upper bound of the MSE of a stationary RW. More precisely, let  $(X_1, \dots, X_n)$  be a sequence of  $n$  vertices sampled by a stationary RW. A RW is stationary iff  $X_1 \sim \pi$ . Let  $(Y_1, \dots, Y_n)$  be a sequence of RE sampled vertices. We show that the MSE, eq.(3), of  $(X_1, \dots, X_n)$  is upper bounded by a function of the MSE of  $(Y_1, \dots, Y_n)$  and  $\alpha$ , where  $0 \leq \alpha < 1$  is the absolute value of the second most dominant eigenvalue of the RW transition probability matrix.

In what follows we define the magnitude of the second most dominant eigenvalue of the RW. Let  $\mathbf{A} = [a_{ij}]$ ,  $i = 1, \dots, |V|$ , be the adjacency matrix of  $G$ ,  $a_{ij} = 1$  iff  $(v_i, v_j) \in E$ , otherwise  $a_{ij} = 0$ . Let

$$\mathbf{D} = \begin{bmatrix} d_{v_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & d_{v_{|V|}} \end{bmatrix}$$

be a diagonal matrix whose diagonal elements are the degrees of the vertices in  $G$ . Let  $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$  be the one-step RW transition probability matrix. The probability that a RW reaches vertex  $v$  from  $u$  in  $t$  steps is

$$p_{uv}^{(t)} = (\mathbf{P}^t)_{uv}.$$

The stationary distribution of the RW is  $\pi = \mathbf{P}\pi$ . Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{|V|}$  be the eigenvalues of  $\mathbf{P}$ . It follows from the fact that  $G$  is an undirected connected non-bipartite graph (and  $\mathbf{P}$  is a stochastic matrix) and the

Frobenius-Perron Theorem that  $1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_{|V|} > -1$  [14]. The absolute value of the second most dominant eigenvalue is defined as

$$\alpha \triangleq \max(\lambda_2, -\lambda_{|V|}). \quad (4)$$

A RW is fast mixing when  $\alpha$  is sufficiently small (we choose to use a vague definition of fast mixing as there are many contradicting definitions of ‘‘fast mixing’’ in the literature).

In the following theorem (Theorem III.1) we show that the estimation error of a RW can be upper bounded by the estimation error of RE sampling and  $\alpha$ .

**Theorem III.1.** *Let  $G = (V, E)$  be an undirected connected non-bipartite graph. Let  $(X_1, \dots, X_n)$  be a sequence of vertices sampled by a stationary RW on  $G$ ,  $n \geq 1$ . Let  $(Y_1, \dots, Y_n)$  be a sequence of RE sampled vertices. Let*

$$\hat{F}(Z_1, \dots, Z_n) \triangleq \frac{1}{n} \sum_{t=1}^n \frac{f(Z_t)}{\pi_{Z_t}}, \quad Z_i \in V, i = 1, \dots, n$$

and let  $\alpha$  be the absolute value of the second most dominant eigenvalue of the RW transition probability matrix.

Then

$$\text{var}(\hat{F}(X_1, \dots, X_n)) \leq \frac{\text{var}(\hat{F}(Y_1, \dots, Y_n))}{(1 - \alpha)}. \quad (5)$$

*Proof:* Let  $\mathbf{S} = \pi^{1/2} \mathbf{P} \pi^{1/2}$ . It is easy to verify that  $\mathbf{S}$  is a  $|V| \times |V|$  symmetric matrix whose eigenvalues are also the eigenvalues of  $\mathbf{P}$ . The eigenvector of  $\mathbf{S}$  corresponding to eigenvalue  $\lambda_1 = 1$  is  $\pi^{1/2}$ . The Courant-Fischer theorem [10, Theorem 4.2.11, pp. 179] gives the second largest eigenvalue of  $\mathbf{S}$

$$\lambda_2 = \max_{w: \langle w, \pi \rangle = 0} \frac{\sum_{\forall v \in V} \sum_{\forall u \in V} w(v)w(u)\pi_v p_{v,u}}{\sum_{\forall u \in V} w(u)^2 \pi_u} \quad (6)$$

and the smallest eigenvalue of  $\mathbf{S}$

$$\lambda_{|V|} = \min_{w: \langle w, \pi \rangle = 0} \frac{\sum_{\forall v \in V} \sum_{\forall u \in V} w(v)w(u)\pi_v p_{v,u}}{\sum_{\forall u \in V} w(u)^2 \pi_u}, \quad (7)$$

as  $\langle w, \pi \rangle = \langle r, \pi^{1/2} \rangle = \sum_{\forall v \in V} r_v \pi_v^{1/2}$ ,  $r = (w(v_i) \sqrt{\pi_{v_i}} : i = 1, \dots, |V|)$ .

Let  $g(v) = f(v)/\pi_v - F$ , which yields  $E[g(X_i)] = 0$ ,  $i = 1, \dots, n$ . In what follows we an upper and lower bound of the covariance of  $g(X_1)$  and  $g(X_t)$ ,  $t = 2, \dots, n$ . Consider the following definitions of covariance and variance: For  $1 < t \leq n$

$$\text{cov}(g(X_1), g(X_t)) \triangleq \sum_{\forall v \in V} \sum_{\forall u \in V} g(v)g(u)\pi_v p_{v,u}^{(t)}$$

and

$$\text{var}(g(X_i)) \triangleq \sum_{\forall u \in V} g(u)^2 \pi_u, \quad i = 1, \dots, n.$$

The bounds are found by replacing the above definitions into eqs.(6) and (7)

$$\lambda_2 \geq \frac{\text{cov}(g(X_1), g(X_2))}{\text{var}(g(X_1))} \geq \lambda_{|V|}.$$

Let  $\alpha = \max(\lambda_2, -\lambda_{|V|})$ , as defined in eq.(4). Then

$$\alpha \geq \frac{\text{cov}(g(X_1), g(X_2))}{\text{var}(g(X_1))}.$$

As  $\lambda_2^t$  and  $\lambda_{|V|}^t$  are eigenvalues of  $P^t$  and  $\text{var}(g(X_1)) > 0$ , we have that

$$\alpha^t \text{var}(g(X_1)) \geq \text{cov}(g(X_1), g(X_t)). \quad (8)$$

A known property of the variance is [19, pp. 265]

$$\begin{aligned} \text{var}\left(\frac{1}{n} \sum_{i=1}^n g(X_i)\right) &= \frac{1}{n} \text{var}(g(X_1)) + \\ &\quad \frac{2}{n} \sum_{t=2}^n \frac{n-t}{n} \text{cov}(g(X_1), g(X_t)). \end{aligned} \quad (9)$$

Applying eq.(8) into eq.(9) yields

$$\begin{aligned} \text{var}\left(\frac{1}{n} \sum_{i=1}^n g(X_i)\right) &\leq \text{var}(g(X_1)) \left(\frac{1}{n} + \frac{2}{n} \sum_{t=2}^n \frac{n-t}{n} \alpha^t\right) \\ &\leq \text{var}(g(X_1)) \left(\frac{1}{n} + \frac{2\alpha}{n(1-\alpha)}\right) = \\ &= \text{var}(g(X_1)) \frac{1+\alpha}{n(1-\alpha)} \\ &\leq \frac{\text{var}(g(X_1))}{n(1-\alpha)}, \end{aligned}$$

as  $0 \leq \alpha < 1$ ,

$$\sum_{t=2}^n \frac{n-t}{n} \alpha^t = \frac{\alpha}{1-\alpha} - \frac{2\alpha^2 - \alpha^3 - \alpha^{n+1}}{n(1-\alpha)^2}$$

and  $2\alpha^2 - \alpha^3 - \alpha^{n+1} \geq 0$ .

The proof is concluded by noting that

$$\frac{1}{n} \text{var}(g(X_1)) = \frac{1}{n} \text{var}(\hat{F}(Y_1)) = \text{var}(\hat{F}(Y_1, \dots, Y_n))$$

and that

$$\text{var}\left(\frac{1}{n} \sum_{i=1}^n g(X_i)\right) = \text{var}(\hat{F}(X_1, \dots, X_n)).$$

The above proof is valid for any value of  $\alpha$ . The upper bound in Theorem III.1 is tight as  $\alpha = \lambda_2 = \lambda_{|V|} = 0$  yields  $\text{cov}(g(X_1), g(X_t)) = 0$ ,  $t > 1$ . Eq.(9) implies  $\text{var}(\hat{F}(X_1, \dots, X_n)) = \text{var}(\hat{F}(Y_1, \dots, Y_n))$ . Hence, the MSE of a fast mixing RW can be approximated by the MSE of RE, but only for small enough values of  $\alpha$ . Otherwise, we need to use the upper bound in Theorem III.1. ■

#### IV. MSE MINIMIZATION

In Section III we provided an upper bound for the RW MSE. So far we considered a RW that samples  $v \in V$  proportional to  $d_v$ , i.e., vertices are sampled from distribution  $\pi$ . In what follows we denote this type of random walk “standard RW”.

There are many different ways to sample a graph (e.g., RE, RV, standard RW). In this section we are interested in sampling the graph as to minimize a given weighted sum of the MSE. The motivation behind this section comes from the several types of stationary RWs that sample vertices,  $(X_i)_{i=1}^n$ , with (an arbitrary) distribution  $X_i \sim \gamma$ ,  $i = 1, \dots, n$  (e.g., Metropolis-Hastings algorithm, Gibbs [19] sampler, and weighted RW are three of such RW types). In Theorem III.1 we proved that the MSE of a standard RW is upper bounded by the MSE of RE sampling times a constant that depends on the graph. Unfortunately, Theorem III.1 cannot be easily extended to  $\gamma \neq \pi$ . Thus, we make the simplifying assumption of independence in  $(X_i)_{i=1}^n$ . For MHRWu the independence assumption means that vertices are RV sampled. For RW the independence assumption means that vertices are RE sampled. In our simulations in Section VII we see that the independence assumption gives a good approximation for RWs (specially if vertices are sampled by our proposed RW, frontier sampling, described in Section VI) and a bad approximation for MHRWu.

To illustrate the optimization, consider estimating the degree distribution,  $\theta_d$ ,  $d = 0, 1, \dots$ ,

$$\theta_d \triangleq \frac{1}{n} \sum_{v \in V} f_d(v),$$

where  $f_d(v) = \mathbf{1}(d_v = d)/|V|$  and  $\mathbf{1}(d_v = d) = 1$  if  $d_v = d$  and zero otherwise, from a sequence of i.i.d. sampled vertices,  $(Y_1, \dots, Y_n)$ , where  $Y_i \sim \gamma$ ,  $i = 1, \dots, n$ . Note that

$$\hat{F}_d(Y_1, \dots, Y_n) = \frac{1}{n} \sum_{t=1}^n \frac{f_d(Y_t)}{\gamma_{Y_t}}.$$

is an unbiased estimate of  $\theta_d$ . To simplify our analysis we assume that the probability of sampling vertex  $v \in V$  only depends on  $d_v$ , the degree of  $v$ , i.e.,  $\gamma_v = \Gamma_{d_v} > 0$ . From the independence of  $Y_i$ ,  $i = 1, \dots, n$ ,

$$\begin{aligned} E[(\hat{F}_d(Y_1, \dots, Y_n) - F_d)^2] &= \frac{1}{n} \left( \sum_{v \in V} \left( \frac{f_d(v)}{\gamma_v} \right)^2 \gamma_v - \theta_d^2 \right) \\ &= \frac{\theta_d}{\Gamma_d} - \theta_d^2. \end{aligned}$$

Let  $\gamma^*$  be the distribution that minimizes the weighted MSE

$$\gamma^* = \arg_{\gamma} \min \sum_{\forall d} \left( \frac{\theta_d}{\Gamma_d} - \theta_d^2 \right) w_d,$$

$w_v > 0$ ,  $v \in V$ .

**Lemma IV.1.** *The distribution  $\Gamma^*$  that minimizes the weighted MSE*

$$\Gamma^* = \arg_{\Gamma} \min \sum_{\forall d} \left( \frac{\theta_d}{\Gamma_d} - \theta_d^2 \right) w_d,$$

with weights  $\{w_j\}$  satisfies the following relation

$$\frac{w_i}{w_j} = \left( \frac{\Gamma_i}{\Gamma_j} \right)^2$$

*Proof:* As  $\Gamma_d$ ,  $d = 1, 2, \dots$ , is a distribution we add the restriction  $\sum_{\forall d} \Gamma_d = 1$  as a Lagrange multiplier in the optimization, which results in the set of equations:

$$h(d) = \sum_{\forall d} \left( \frac{\theta_d}{\Gamma_d} - \theta_d^2 \right) w_d - \lambda \left( \sum_{\forall d} \Gamma_d - 1 \right), \forall u \in V$$

Taking the derivative of  $h(d)$  with respect to  $\Gamma_d$  and equating to zero yields

$$\frac{\partial h(d)}{\partial \Gamma_d} = -\frac{1}{\Gamma_d^2} - \lambda = 0, \forall d,$$

thus

$$\frac{w_i}{w_j} = \left( \frac{\Gamma_i}{\Gamma_j} \right)^2$$

All is left is to prove that  $\frac{w_i}{w_j} = \left( \frac{\Gamma_i}{\Gamma_j} \right)^2$  is not a saddle point of  $h(d)$ . This is easy as  $\partial^2 h(d) / \partial^2 \Gamma_d = 2 / \Gamma_d^3 > 0$ . ■

In particular,  $\Gamma_d \propto d$  gives

$$w_i = w_j (i/j)^2, \quad i, j = 1, 2, \dots$$

Note that  $w_i > w_j$  when  $i > j$  and  $w_i < w_j$  when  $d_i < d_j$ . Thus, a standard RW and RE sampling optimize a weighted MSE that places larger weights at the tail of the degree distribution. Another particular case occurs when  $\Gamma_i = \Gamma_j$ ,  $\forall i, j$ :

$$w_v = w_u, \quad \forall u, v \in V,$$

i.e., RV sampling optimizes a weighted MSE with equal weights.

*Degree distribution: RW v.s. RV sampling*

The NMSE is the Normalized Mean Square Error, defined as  $\sqrt{MSE}/F$ , where  $F$  is the true value. Let  $\bar{d}$  be the average degree. From the exposition in Section IV it is straightforward to show that the NMSE estimating  $\theta_d$  using  $n$  RE samples is

$$\text{NMSE}_{\text{re}}(d) = \sqrt{(\bar{d}/(d\theta_d) - 1)/n}, \quad d > 0. \quad (10)$$

Similarly, the NMSE( $d$ ) using RV sampling is

$$\text{NMSE}_{\text{rv}}(d) = \sqrt{(1/\theta_d - 1)/n}. \quad (11)$$

Applying Theorem III.1 to eq.(10) yields

$$\text{NMSE}_{\text{RW}}(d) \leq \sqrt{\frac{(\bar{d}/(d\theta_d) - 1)}{n(1 - \alpha)}}, \quad d > 0. \quad (12)$$

From equations (12) and (11) we see that a fast mixing RW more accurately estimates degrees larger than the average ( $d > \bar{d}$ ) while RV sampling more accurately estimates degrees smaller than the average ( $d < \bar{d}$ ). The above analysis explains what has been previously observed in [16].

## V. MHRWu v.s. RWs

The RW described in Section III is the most common type of RW found in the literature [14] but there are other types of random walks. For more details refer to [19, Chapter 7]. In [21] a Metropolis-Hastings RW that samples vertices uniformly at random is described, which we denote MHRWu in this work. MHRWu is found to be less accurate than a RW in estimating some graph characteristics such as the degree distribution [7].

A MHRWu is an accept-reject sampling process that samples vertices uniformly. In this section we explore a parallel between MHRWu and a RE resampling algorithm (presented in Section V-A). The MHRWu works as follows, starting at vertex  $v$  we:

- select a neighbor  $u$  of  $v$  uniformly at random;
- the next sampled vertex (step) is  $u$  with probability  $\min(d_v/d_u, 1)$ , otherwise  $v$  is the next (step) sampled vertex. This is equivalent to say that we add a copy of vertex  $v$  to the sample set with probability  $\max(0, 1 - d_v/d_u)$ , otherwise  $u$  is the next (step) sampled vertex.

### A. RW alternative unbiased estimator by resampling

In what follows we present another way to obtain an unbiased estimate of  $F$ . This estimator will be later used to provide an approximation to MHRWu MSE. As before,  $(Y_i)_{i=1}^n$  is a sequence of vertices sampled by RE. Let  $K_d$  be the number of vertices with degree  $d$  in  $(Y_i)_{i=1}^n$ .  $K_d$  is a Binomial random variable with parameters  $n$  and  $p_d = d/\text{vol}(V)$  ( $P[K_d = k] = \binom{n}{k} p_d^k (1 - p_d)^{n-k}$ ). Let  $Z_i^{(d)} \in \{1, 2, \dots\}$  be a sequence of i.i.d. Geometric random variable with parameter  $p_d$ ,  $i = 1, \dots, n$ .

#### Lemma V.1.

$$F' = 1/n \sum_{i=1}^n f(Y_i) Z_i^{(d_{Y_i})},$$

where  $d_{Y_i}$  is the degree of  $Y_i$ , is an unbiased estimate of  $F$  (eq.(1)).

*Proof:* First

$$\begin{aligned} E[f(Y_i) Z_i^{(d_{Y_i})}] &= E[E[f(Y_i) Z_i^{(d_{Y_i})} | Y_i]] = \\ &= E[E[f(Y_i) | Y_i] E[Z_i^{(d_{Y_i})} | Y_i]] = \\ &= E[f(Y_i) E[Z_i^{(d_{Y_i})} | Y_i]] = E[f(Y_i) 1/\pi_{Y_i}] = \\ &= \sum_{v \in V} f(v) \pi_v / \pi_v = \sum_{v \in V} f(v). \end{aligned}$$

As  $(Y_i)_{i=1}^n$  is an i.i.d. sequence of random variables  $E[F'] = (1/n)nE[f(Y_i) Z_i^{(d_{Y_i})}]$ , which concludes our proof. ■

The next lemma provides the NMSE of estimating  $\theta_j$  using Lemma V.1.

#### Lemma V.2.

$$\text{NMSE}_{F'}(d) = \sqrt{\frac{2(1 - p_d)}{np_d \theta_d^2}}.$$

*Proof:* The total number of replications of vertices with degree  $d$  is  $Z_d = \sum_{j=1}^{K_d} Z_j^{(d)}$ . Note that  $F' = Z_d/n$ , which yields [20, pp. 349, Example 4n]

$$\begin{aligned} \text{var}(F') &= \frac{1}{n^2} \text{var}\left(\sum_{j=1}^{K_d} Z_j^{(d)}\right) = \\ &= (1/n^2) \left(E[K_d] \text{var}(Z_j^{(d)}) + E[Z_j^{(d)}]^2 \text{var}(K_d)\right) = \\ &= (1/n^2) \left((np_d)(1 - p_d)/p_d^2 + (1/p_d^2)(np_d(1 - p_d))\right) = \\ &= (1/n^2) \left(n(1 - p_d)/p_d + n(1 - p_d)/p_d\right) = \frac{2(1 - p_d)}{np_d}. \end{aligned}$$

Lemma V.1 gives  $E[F'] = \theta_d$ , which yields

$$\text{NMSE}_{F'}(d) = \sqrt{\text{var}(F')}/\theta_d = \sqrt{\frac{2(1 - p_d)}{np_d \theta_d^2}}. \quad \blacksquare$$

### B. Metropolis-Hasting RW: Uniform vertex sampling (MHRWu)

We now turn our attention to MHRWu of [21]. A different way to present the MHRWu algorithm is as an edge sampling process that samples vertices uniformly. The MHRWu is time reversible, which means that the same number of walkers going from  $v$  to  $u$  must go from  $u$  to  $v$ . Let  $p_{ab}$  be the probability that the walker goes from  $a$  to  $b$ ,  $a, b \in V$ . Vertices are sampled uniformly and the Markov chain is time reversible, which yields  $p_{vu} = p_{uv}$ . To simplify our exposition we assume, without loss of generality, that  $d_u < d_v$ . It is easy to see that  $p_{vu} = p_{uv} = 1/d_v$  satisfies the above conditions and can be implemented by selecting neighbors of  $v$  and  $u$

uniformly at random. However, because the walker chooses  $v$  with probability  $1/d_u > 1/d_v$  we are required to add a self-loop with probability  $1/d_u - 1/d_v$  at  $u$ , as illustrated in Figure 1(a). The arrows in Figure 1(a) indicate the walker direction and the probability that the direction is taken. Thus, the probability that an edge  $(v, u) \in E$  is sampled is  $1/(d_v|V|)$ . The self-loop adds an average of  $(d_v - d_u)/(d_v(d_u - 1) + d_u)$  “extra” copies of  $u$  for each sample of  $v$ , due to the existence of the edge  $(u, v)$ .

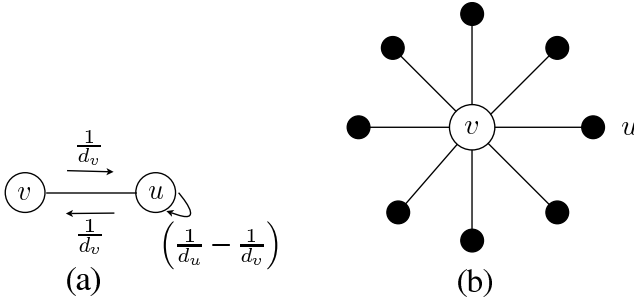


Fig. 1. (a) MHRWu transition probabilities and (b) star graph example.

To illustrate the problem with MHRWu consider the graph in Figure 1(b). The self-loop at  $u$  has probability  $1 - 1/d_v$ , which means that on average  $d_v - 1$  extra copies of  $u$  are made for each sample of  $v$ . Note that the resampling of vertex  $u$  can be described by the resampling algorithm in Section V-A. However, as seen next, in general a MHRWu resamples vertices significantly less often than the resampling algorithm described in Section V-A.

### C. A MHRWu MSE approximation

Consider an edge-sampling process that samples edge  $(u, v)$  with probability  $1/(\max(d_v, d_u)|V|)$ ,  $\forall (u, v) \in E$ . Edges are sampled independently. After edge  $(u, v)$  is selected, the probability that node  $v$  is resampled is  $1/d_v - 1/d_u$  if  $d_v < d_u$  and zero otherwise. The same is valid for node  $u$ . While in a MHRWu only the vertex with a self-loop is allowed to *resample* (i.e., make multiple copies of itself), we simplify our analysis by assuming that edge  $(u, v)$  incur self-loops on both  $u$  and  $v$  with probabilities  $1/d_u$  and  $1/d_v$ , respectively. Let  $(Y_i)_{i=1}^n$  be a sequence of vertices sampled by RE and let  $K_d$  be the number of vertices with degree  $d$  in  $(Y_i)_{i=1}^n$ . Setting  $p_d = 1/d$  in

Lemma V.2 we get

$$\begin{aligned} \text{NMSE}'_{\text{mh}}(d) &\approx \sqrt{\frac{E[K_d]\text{var}(Y_j^{(d)}) + E[Y_j^{(d)}]^2\text{var}(K_d)}{n^2\theta_d^2}} = \\ &\sqrt{\frac{n\theta_d(d^2 - d) + d^2n\theta_d(1 - \theta_d)}{n^2\theta_d^2}} = \sqrt{\frac{d^2(2 - \theta_d) - d}{n\theta_d}} \\ &> \sqrt{\frac{2(d-1)^2/\theta_d - 1}{n}}. \end{aligned}$$

Although there are no guarantees that  $\text{NMSE}'_{\text{mh}}(d)$  is a good approximation to the true NMSE of MHRWu, our simulations (some of these results presented in Section VII) show that  $\text{NMSE}'_{\text{mh}}(d)$  is indeed close to the empirical value of NMSE for large enough values of  $d$ .

It is interesting to take a closer look at the equation

$$\text{NMSE}'_{\text{mh}}(d) > \sqrt{\frac{(d-1)^2/\theta_d - 1}{n}}$$

to note that the MHRWu NMSE grows linearly with  $d$ . Contrast the linear growth of  $\text{NMSE}'_{\text{mh}}(d)$  as a function of  $d$  with the NMSE of a RW:

$$\text{NMSE}_{\text{rw}}(d) \leq \sqrt{(\hat{d}/(d\theta_d) - 1)/(n(1 - \alpha))},$$

which is inversely proportional to  $d$ . Also note that

$$\text{NMSE}'_{\text{mh}}(d) > \text{NMSE}_{\text{rw}}(d), \forall d,$$

in fact, the error of MHRWu is almost  $d$  times larger than the error of RW sampling.

There is still room for improvement in our approximation. MHRWu is a type of random walk and that the NMSE of MHRWu, similar to the NMSE of RW, should also increase with  $\alpha$ . Our approximation considers independently sampled edges (i.e.,  $\alpha = 0$ ). As part of future work we intend to derive the MSE expression of MHRWu that includes  $\alpha$ . Table I summarizes our analytical results.

## VI. CURRENT EFFORTS TO IMPROVE RW ACCURACY

Sampling a graph using a RW is not without drawbacks. A random walker can get (temporarily) “trapped” inside a subgraph whose characteristics differ from those of the whole graph. Even if the random walker starts in steady state (i.e., is stationary), a “trap” may increase the mean squared error of the estimates. Ideally, the random walker needs to mitigate the effect of these traps on the estimates. Note that in a graph with such “traps”  $\alpha \approx 1$  and, as seen in Section III, the RW MSE is upper bounded by  $1/(1 - \alpha)$ . A simple naive solution to the RW “trapping” problem (adopted in [7] to sample Facebook), is to sample the graph using multiple independent random walkers [5]. This naive solution, however, can have the opposite effect and exacerbate the problem [18]. The literature, however,

Sampling Method	NMSE error
Node Sampling	$\sqrt{\frac{1/\theta_d-1}{n}}$
Edge Sampling	$\sqrt{\frac{\bar{d}/(d\theta_d)-1}{n}}, d > 0$
Random Walk	$\leq \sqrt{\frac{(\bar{d}/(d\theta_d)-1)}{n(1-\alpha)}}, d > 0$
Metropolis-Hastings RW	$> \sqrt{\frac{(d-1)^2/\theta_d-1}{n}}, d > 0$

TABLE I

SUMMARY OF RESULTS: DEGREE DISTRIBUTION ESTIMATION ERRORS OF VARIOUS SAMPLING METHODS.  $\theta_d$  IS THE FRACTION OF NODES WITH DEGREE  $d$  (QUANTITY THAT IS ESTIMATED),  $n$  IS THE NUMBER OF SAMPLED NODES, AND  $\alpha$  THE ABSOLUTE VALUE OF THE SECOND MOST DOMINANT EIGENVALUE OF THE RW TRANSITION PROBABILITY MATRIX

provides some promising approaches to cope with this problem if the graph admits a limited (small) amount of RV sampling.

RV sampling has been used to significantly reduce  $\alpha$  [1] by allowing the random walks to “jump” to an RV sampled node. The algorithm in [1] differs from the PageRank [3] RW + RV “jumps” in that it obtains unbiased estimates of eq.(1). In [1] it is also shown that, unless the underlying graph is known, the PageRank algorithm must necessarily obtain biased estimates of eq.(1). Another promising approach to improving the RW accuracy is starting  $m$  dependent walkers at  $m$  RV sampled nodes. This approach, called Frontier Sampling (FS) [18] given in Algorithm 1, introduces a simple dependence among all  $m$  walkers in a way that starting the  $m$  walkers at  $m$  RV sampled nodes is arbitrarily close to starting FS in steady state, provided  $m$  is large enough. In our simulations we observe that the FS NMSE for  $m = 1000$  is close to the NMSE of a RW with negligible mixing time.

---

**Algorithm 1:** Frontier Sampling (FS).

---

- 1:  $n \leftarrow 0$  { $n$  is the number of steps}
  - 2: Initialize  $L = (v_1, \dots, v_m)$  with  $m$  randomly chosen vertices (uniformly)
  - 3: **repeat**
  - 4:   Select  $u \in L$  with probability  $d_u / \sum_{v \in L} d_v$
  - 5:   Select an outgoing edge of  $u$ ,  $(u, v)$ , uniformly at random
  - 6:   Replace  $u$  by  $v$  in  $L$  and add  $(u, v)$  to sequence of sampled edges
  - 7:    $n \leftarrow n + 1$
  - 8: **until**  $n \geq n - mc$
- 

## VII. SIMULATION RESULTS

In what follows we present the results of our simulations of the sampling methods discussed in this paper. The graphs used in our experiments are real-world graphs detailed in Table II. But due to space constraints we restrict our results to the two largest graphs in our datasets: LiveJournal and Flickr. Note that our simulations are performed on disconnected graphs, which can increase the MSE of methods such as RW and MHRWu (FS is designed to mitigate the large MSEs caused by disconnected graphs). The results using the other datasets are similar to LiveJournal and Flickr results, no surprises worth reporting. All sampling methods have a budget of  $n$  vertices to sample. Each newly sampled vertex deducts one from the budget while resampling a vertex does not change the budget (i.e., has cost zero). The empirical MSE of our simulations is obtained over 10,000 runs.

In some of our simulations we use a slightly different MSE metric than the NMSE: the normalized root mean square error of the Complementary Cumulative Distribution Function (CCDF)  $\gamma = \{\gamma_d\}_{d \geq 1}$ , where  $\gamma_d = \sum_{k=d+1}^{\infty} \theta_k$ ,

$$\text{CNMSE}(d) = \frac{\sqrt{E[(\hat{\gamma}_d - \gamma_d)^2]}}{\gamma_d}, \quad (13)$$

where  $\hat{\gamma}_d$  is the estimate of  $\gamma_d$ . The CNMSE is just the NMSE of  $\gamma_d$  and thus  $\text{CNMSE}'_{\text{mh}}(d)$ ,  $\text{CNMSE}_{\text{rv}}(d)$ ,  $\text{CNMSE}_{\text{rw}}(d)$ , and  $\text{CNMSE}_{\text{re}}(d)$  have the same equation as their respective NMSE formulation with  $\theta_d$  replaced by  $\gamma_d$  (or  $\Pi_d = d\theta_d/\bar{d}$  replaced by  $d\gamma_d/\bar{\gamma}$ ).

Note that the graphs in Table II are directed. Obtaining directed graph characteristics such as the in-degree distribution from graphs that can be crawled like undirected graphs (e.g., Twitter and Livejournal) is a trivial task, for more details refer to [18].

### Goodness of theoretical approximations

*FS v.s. RW and RV sampling:* This first set of simulations differ from the remaining simulations in this paper in that resampling a vertex reduces the sampling budget by one. In our results we compare the MSE of FS and RE. In our first result, Figure 2 shows the log-log plot of the in-degree NMSE of FS, RW, and RV sampling. In Figure 2 we observe that the FS NMSE is close to the RE NMSE for all degrees  $d > 0$  (note that from Theorem III.1 the RE NMSE is equivalent to the MSE of a RW with negligible mixing time ( $\alpha \ll 1$ )). The same is true in all other datasets. Moreover, as theoretically predicted by the analysis performed in Section IV, the NMSE of RV is smaller than the NMSE of RE when  $d$  is smaller than the average degree and the NMSE of RV is larger than the NMSE of RE when  $d$  is larger than the average degree.

Graph	Flickr	LiveJournal	YouTube	Internet RLT
Description	Social Net.	Social Net.	Social Net.	Internet tracent.
Type of graph	Directed	Directed	Directed	Directed
# of Vertices	1,715,255	5,204,176	1,138,499	192,244
Size of LCC	1,624,992	5,189,809	1,134,890	190,914
# of Edges	22,613,981	77,402,652	9,890,764	609,066
Average Degree	12.2	14.6	8.7	3.2
$w_{\max}$	2232	1029	3305	335
% of Original Graph	26.9%	95.4%	NA	NA

TABLE II

SUMMARY OF THE GRAPH DATASETS USED IN OUR SIMULATIONS. “SIZE OF LCC” REFERS TO THE SIZE OF THE LARGEST CONNECTED COMPONENT AND  $w_{\max}$  IS THE VALUE OF THE LARGEST VERTEX DEGREE DIVIDED BY THE AVERAGE DEGREE.

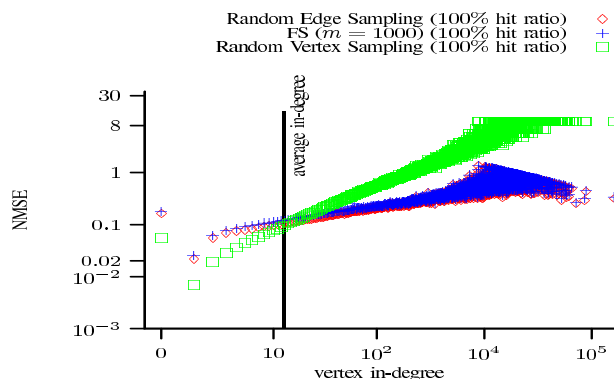


Fig. 2. (Flickr) The log-log plot shows the NMSE of the in-degree distribution estimation with budget  $n = |V|/100 = 18612$  (NMSE over 10,000 runs).

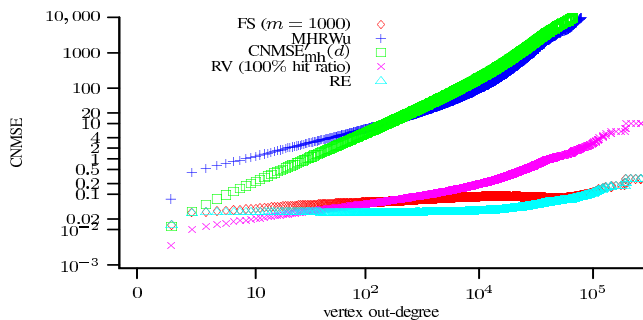


Fig. 3. (Flickr) The log-log plot of the CNMSE of the in-degree distribution estimates with budget  $n = |V|/100$ . RV with hit ratio 100%.

*FS v.s.  $CNMSE_{rw}(d)$  bound & MHRWu v.s.  $NMSE'_{mh}(d)$  and RV sampling:* In this simulation on Flickr we seek to assess the goodness of the theoretical approximation of  $NMSE'_{mh}(d)$  derived in Section V-C. We also seek to test if assuming that the CNMSE of FS is equivalent to the CNMSE of a RW that mixes fast (i.e.,  $\alpha$  is small). We simulate FS, MHRWu, and RV on the LiveJournal graph with  $n = |V|/100$  samples each. Figure 3 plots the the in-degree CNMSE of FS, MHRWu, RV and also plots  $CNMSE'_{re}(d)$  and  $CNMSE'_{mh}(d)$ . Note that the approximation of  $CNMSE'_{mh}(d)$  is accurate

when  $d \geq 10^2$ . Remark that the CNMSE of MHRWu is so large that for vertices with degree greater than  $8 \times 10^4$  it is greater than 10,000. The estimates of MHRWu are clearly much less accurate than the estimates of FS. As we do not have  $\alpha$ , we consider  $\alpha = 0$ , i.e.,  $CNMSE_{rw}(d) \approx CNMSE_{re}(d)$ . Note that  $CNMSE_{re}(d)$  approximates well the CNMSE of FS. From Theorem III.1 we know that this is equivalent to the CNMSE of a RW with  $\alpha$  small. Also observe that the CNMSE of MHRWu is much larger than the CNMSE of RV and therefore, as expected, the RV CNMSE is not a good approximation to the MHRWu CNMSE.

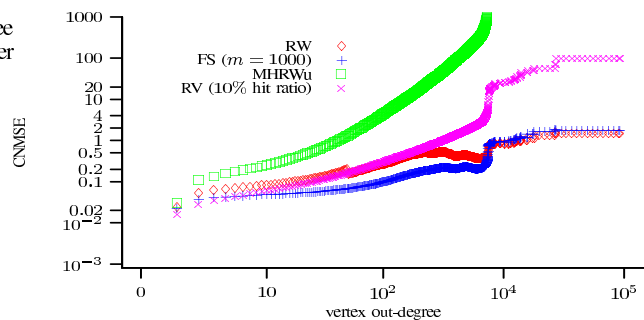


Fig. 4. (LiveJournal) The log-log plot of the CNMSE of the in-degree distribution estimates with budget  $n = |V|/1000$ .

### Accuracy of Graph Sampling Methods

In this simulation we compare RW, FS, MHRWu, and RV (with 10% hit ratio). Figure 4 plots the in-degree CNMSE of RW, FS, MHRWu, and RV (with 10% hit ratio) on the LiveJournal graph for budget of  $n = |V|/1000$ . “RV with (with 10% hit ratio)” represents random vertex sampling when only 1 in 10 queries are valid, i.e., in average only  $n/10$  samples are used in the estimator. We observe that RV (with 10% hit ratio) is less accurate than RW and FS. FS is slightly more accurate than RW for degrees between 10 and  $5 \times 10^3$ .

A similar simulation with  $n = |V|/100$  on the Flickr graph reveals a similar picture (results shown in Figure 5).



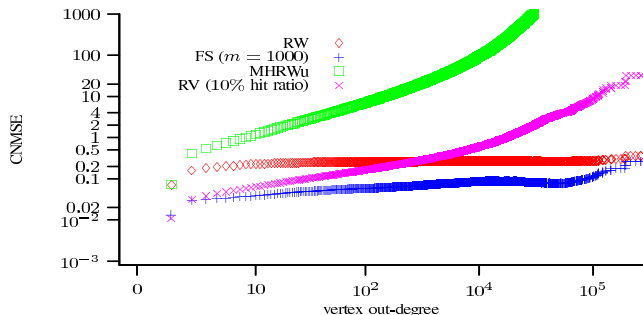


Fig. 5. (Flickr) The log-log plot of the CNMSE of the in-degree distribution estimates with budget  $n = |V|/100$ .

RV (with 10% hit ratio) is the most accurate sampling method for degree  $d = 1$  (with FS in a close second place). For degrees  $d > 1$  FS is the most accurate method. RW, however, performs poorly when compared to FS (the CNMSE is up to one order of magnitude larger). MHRWu is again the least accurate method (where  $\text{CNMSE}(d) > 1000$  when  $d > 2 \times 10^4$ ).

### VIII. RELATED WORK

The first work to use resampling of RE as a rough approximation to MHRWu sampling was [6]. We, however, refine this rough approximation by computing the exact probability of a self-loop, obtaining a good estimate of the average number of resamples of the same node.

### IX. CONCLUSIONS

This paper provides an upper bound for the MSE of a stationary RW as a function of the MSE of RE and the absolute value of the second most dominant eigenvalue of the RW transition probability matrix. We observed that RW and RV sampling are optimal in respect to different weighted MSE optimizations and analyzed when RW is preferable to RV sampling. We also presented an approximation to the MHRWu MSE. Finally, we introduce a novel RW sampling algorithm, Frontier Sampling (FS). Our simulation experiments on large real world graphs showed that FS achieves the MSE of a RW with negligible mixing time.

### X. ACKNOWLEDGMENTS

We would like to thank Maciej Kurant for the helpful discussions. We also would like to thank Ananthram Swami for helpful discussions regarding the MSE optimization and Alan Mislove for kindly making available some of the data used in this paper.

### REFERENCES

- [1] Konstantin Avrachenkov, Bruno Ribeiro, and Don Towsley. Improving random walk estimation accuracy with uniform restarts. In *Proc. of the 7th Workshop on Algorithms and Models for the Web Graph*, 2010.
- [2] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang. Complex networks: Structure and dynamics. *Physics Reports*, 424(4-5):175–308, 2006.
- [3] S. Brin and L. Page. The anatomy of a large-scale hypertextual Web search engine. In *Proc. of the WWW*, 1998.
- [4] Nathan Eagle, Alex S. Pentland, and David Lazer. Inferring friendship network structure by using mobile phone data. *PNAS*, 106(36):15274–15278, August 2009.
- [5] Charles J. Geyer. Practical Markov Chain Monte Carlo. *Statistical Science*, 7(4):473–483, 1992.
- [6] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Practical recommendations on sampling on users by crawling the social graph. *JSAC special issue on Measurement of Internet Topologies*, 2011.
- [7] Minas Gjoka, Maciej Kurant, Carter T. Butts, and Athina Markopoulou. A walk in Facebook: Uniform sampling of users in online social networks. In *Proc. of the IEEE Infocom*, March 2010.
- [8] Christos Gkantsidis, Milena Mihail, and Amin Saberi. Random walks in peer-to-peer networks: algorithms and evaluation. *Perform. Eval.*, 63(3):241–263, March 2006.
- [9] Monika R. Henzinger, Allan Heydon, Michael Mitzenmacher, and Marc Najork. On near-uniform url sampling. In *Proceedings of the WWW*, pages 295–308, 2000.
- [10] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 1990.
- [11] Marlon A. Konrath, Marinho P. Barcellos, and Rodrigo B. Mansilha. Attacking a swarm with a band of liars: evaluating the impact of attacks on bittorrent. In *P2P '07: Proceedings of the Seventh IEEE International Conference on Peer-to-Peer Computing*, pages 37–44, Washington, DC, USA, 2007. IEEE Computer Society.
- [12] Jure Leskovec and Christos Faloutsos. Sampling from large graphs. In *Proc. of the KDD*, pages 631–636, 2006.
- [13] Jure Leskovec, Kevin J. Lang, Anirban Dasgupta, and Michael W. Mahoney. Statistical properties of community structure in large social and information networks. In *Proc. of the WWW*, pages 695–704, 2008.
- [14] L. Lovász. Random walks on graphs: a survey. *Combinatorics*, 2:1–46, 1993.
- [15] Alan Mislove, Massimiliano Marcon, Krishna P. Gummadi, Peter Druschel, and Bobby Bhattacharjee. Measurement and Analysis of Online Social Networks. In *Proc. of the IMC*, October 2007.
- [16] Amir H. Rasti, Mojtaba Torkjazi, Reza Rejaie, Nick Duffield, Walter Willinger, and Daniel Stutzbach. Respondent-driven sampling for characterizing unstructured overlays. In *Proc. of the IEEE Infocom*, pages 2701–2705, April 2009.
- [17] Bruno Ribeiro, William Gauvin, Benyuan Liu, and Don Towsley. On MySpace account spans and double Pareto-like distribution of friends. In *IEEE Infocom 2010 Network Science Workshop*, Mar 2010.
- [18] Bruno Ribeiro and Don Towsley. Estimating and sampling graphs with multidimensional random walks. In *Proc. of the ACM SIGCOMM IMC*, Oct. 2010.
- [19] Christian P. Robert and George Casella. *Monte Carlo Statistical Methods*. Springer-Verlag, 2nd edition, 2005.
- [20] Keith Ross. *A First Course in Probability*. Prentice Hall, 5 edition, 1997.
- [21] Daniel Stutzbach, Reza Rejaie, Nick Duffield, Subhabrata Sen, and Walter Willinger. On unbiased sampling for unstructured peer-to-peer networks. *IEEE/ACM Trans. Netw.*, 17(2):377–390, 2009.
- [22] Erik Volz and Douglas D. Heckathorn. Probability based estimation theory for Respondent-Driven Sampling. *Journal of Official Statistics*, 2008.