

# Albatross Sampling: Robust and Effective Hybrid Vertex Sampling for Social Graphs

Long Jin<sup>1</sup>, Yang Chen<sup>2</sup>, Pan Hui<sup>3</sup>, Cong Ding<sup>2</sup>, Tianyi Wang<sup>1</sup>,

Athanasios V. Vasilakos<sup>4</sup>, Beixing Deng<sup>1</sup> and Xing Li<sup>1</sup>

<sup>1</sup> Department of Electronic Engineering, Tsinghua University, Beijing, China

<sup>2</sup> Institute of Computer Science, University of Goettingen, Goettingen, Germany

<sup>3</sup> Deutsche Telekom Laboratories/TU-Berlin, Germany

<sup>4</sup> University of Western Macedonia, Greece

jjin-l06@mails.tsinghua.edu.cn, yang.chen@cs.uni-goettingen.de

## ABSTRACT

Nowadays, Online Social Networks (OSNs) have become dramatically popular and the study of social graphs attracts the interests of a large number of researchers. One critical challenge is the huge size of the social graph, which makes the graph analyzing or even the data crawling incredibly time consuming, and sometimes impossible to be completed. Thus, graph sampling algorithms have been introduced to obtain a smaller subgraph which reflects the properties of the original graph well. Breadth-First Sampling (BFS) is widely used in graph sampling, but it is biased towards high-degree vertices during the process of sampling. Besides, Metropolis-Hasting Random Walk (MHRW), which is proposed to get unbiased samples of the social graph, requires the graph to be well connected. In this paper, we propose a vertex sampling algorithm, so-called Albatross Sampling (AS), which introduces random jump strategy into MHRW during the sampling process. The embedded random jump makes the sampling procedure more flexible and avoids being trapped in some locally well connected part. According to our evaluation, we find that no matter using tightly or loosely connected graphs, AS performs significantly better than MHRW and BFS. On the one hand, AS estimates the degree distribution with much lower Normalized Mean Square Error (NMSE) by consuming the same resource budget. On the other hand, to get an acceptable estimation of the degree distribution, AS requires much less resource budget.

## Categories and Subject Descriptors

G.3 [Probability and Statistics]: Statistical computing

## General Terms

Algorithms, Measurement

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

HotPlanet'11, June 28, 2011, Bethesda, Maryland, USA.

Copyright 2011 ACM 978-1-4503-0742-0/11/06 ...\$10.00.

## Keywords

Graph Sampling, Online Social Networks, Random Walk, Metropolis-Hasting Algorithm, Jump Strategy, Degree Distribution, Convergence Time

## 1. INTRODUCTION

In recent years, online social networks (OSNs) have become dramatically popular all over the world. For example, Twitter, which provides microblogging service, has more than 200 million users around the world and generates 110 million “tweets” per day by January, 2011 [7]. The spreading of OSNs also attracts a large number of researchers to explore and study the newly-built large scale networks and their research topics are diverse and fascinating, such as social interaction [23], information propagation [11] and user behavior characterization [4] in social networks.

Social networks are usually modeled as social graphs for analysis. One of the challenges that researchers face is the huge size of OSNs. Firstly, it is almost impractical to deal with the complete datasets, since crawling large social graphs is incredibly time consuming, and sometimes is impossible to be completed. Moreover, algorithms for analyzing huge social graphs require a large quantity of time to compute, even if running on the platform of high-performance computer clusters. Secondly, complete datasets from OSNs are usually not publicly accessible, due to privacy settings of users and the protection from OSN companies. Finally, the number of users in OSNs grows rapidly and connections between users vary with time, thus dynamic large graphs cannot be crawled totally. Therefore, attention has focused on how to shrink a huge social graph to a representative sample, which should have a relatively small size and maintain properties of the original social graph.

Several graph sampling methods have been proposed to obtain samples of OSNs and these graph sampling methods are classified as edge sampling method and vertex sampling method [12]. Random Vertex Sampling and Random Edge Sampling are two basic and intuitive sampling methods. As their names imply, Random Vertex Sampling and Random Edge Sampling methods sample vertices and edges randomly in the whole graph, respectively. However, these two sampling methods are either resource intensive or impractical to sample OSNs [19].

One practical graph sampling method is Breadth-First Sampling (BFS), which is used for social network analysis

in [2, 11, 16, 23]. However, BFS is biased towards high-degree vertices [9]. Another widely used sampling algorithm is Random Walk (RW) and many practical graph sampling methods are based on RW.

On the one hand, RW naturally samples edge uniformly in non-bipartite undirected graphs, but not in disconnected graphs [14]. Frontier Sampling (FS) [19], which is an edge sampling method using multidimensional random walkers, is proposed to exhibit lower estimation errors than RW in the presence of disconnected or loosely connected subgraphs. Since FS is an edge sampling method, it is more convenient to estimate edge-centric properties and specific estimators of degree distribution and global clustering coefficient are proposed in [19].

On the other hand, RW method is biased towards high-degree vertices [9] and Metropolis-Hasting Random Walk (MHRW) method is proposed to obtain unbiased samples of social graphs. MHRW is a vertex sampling methods, whose goal is to mimic random vertex sampling by random walking in graphs. However, one design assumption of MHRW is that the social graph is well connected [8], which results in MHRW not fit for sampling disconnected or loosely connected graphs. Furthermore, Random Jump (RJ), which may jump to any random vertex with a fixed probability in each step, is proposed in [12] to get rid of the walker being stuck in some locally well connected part of the graph.

In this paper, we focus on proposing an improved vertex sampling method which performs well in sampling social graphs, either tightly or loosely connected. By introducing random jump strategy into MHRW, we propose a hybrid sampling algorithm named *Albatross Sampling* (AS). According to our evaluation, given the same sampling cost, AS is more robust for estimating degree distribution with lower Normalized Mean Square Error (NMSE) and also more effective for converging more quickly with much smaller convergence time. Therefore, AS is a promising, robust and effective hybrid vertex sampling algorithm for social network analysis.

This paper is organized as follows. In Section 2, basic properties about social graphs and measurement of information retrieval cost in the sampling process are introduced. Then, in Section 3, existing graph sampling algorithms are discussed in detail and our improved algorithm Albatross Sampling is proposed in Section 4. Finally, the performance of these algorithms for sampling social graphs is evaluated in Section 5 and conclusion is made in Section 6.

## 2. BACKGROUND

### 2.1 Properties of Social Graphs

OSNs are usually modeled as social graphs, which is represented as  $G = (V, E)$ . Here, a vertex  $v$  in set  $V$  represents a user in the OSN and an edge  $e$  in set  $E$  represents a “following” relationship or a friendship link between users, which can be either directed or undirected.

For directed graphs, we define  $k_{iv}$  as the in-degree and  $k_{ov}$  as the out-degree of vertex  $v$ , and  $k_v$  as the degree of vertex  $v$  when the directed graph is changed into undirected graph. Then, we define  $\theta_{ik}$  and  $\theta_{ok}$  as the fraction of vertices with in-degree and out-degree, respectively, less than or equal to  $k$ . Then,  $\widehat{\theta_{ik}}$  and  $\widehat{\theta_{ok}}$  are the corresponding estimated value through sampling.

In the real world, some large scale social network graph-

s are not fully connected and may contain disconnected or loosely connected components, e.g. wireless social networks [6]. The simple random walker may be trapped in some locally well connected part of the graph, and if the properties of that part differ significantly from those of the whole graph, the sample set cannot represent the original graph well. Therefore, robust and effective methods for sampling disconnected or loosely connected graphs should be studied.

### 2.2 Measurement of Sampling Cost

In this part, we propose several definitions related to the measurement of information retrieval cost during sampling.

Firstly, Total-Cost is defined as the total resource budget (time, bandwidth, or cache) that is used in the process of sampling. Since downloading the profile of a user is much more time-consuming compared to making the choice of the next sampled user, Total-Cost can be viewed as the maximum number of unique users that are visited during sampling.

Secondly, in most social networks, all incoming and outgoing neighbors of a sampled vertex can be learned. For instance, when we visit a user in Twitter, all the followers and followees in the user profile can be known. This fact makes graph sampling methods, such as BFS and RW, feasible and cheap. We define the resource for visiting a new neighbor of the current vertex as Walk-Cost, which can be normalized as 1. We should mention that information of sampled vertices are stored in cache, therefore only visiting new vertices generates sampling cost.

Finally, in different sampling methods, various sampling strategies may be adopted in each step and more resource may be used. For example, in MySpace and Twitter, each user is given a unique user-ID, thus Random Vertex Sampling can be implemented by selecting a random number from the space of user-IDs in each step. If the selected number is a valid user-ID, the user is sampled, otherwise the number is discarded. However, for some OSNs, valid user-IDs are sparse in the space of user-IDs. Sampling a valid user requires nearly 10 attempts on average in myspace [18] while only 1.5 attempts on average in Twitter according to the measurement results provided in [11]. Similarly, in Random Jump, we may jump to any vertex to restart the random walker with a fixed probability in each step. To compare different sampling algorithms fairly, we define the resource for choosing a new random vertex in the total graph as Jump-Cost and we choose Jump-Cost as 10 in our evaluation.

## 3. ANALYSIS OF EXISTING SAMPLING ALGORITHMS

The sampling process of social graphs usually starts from either one or several initial vertices, which can be called seeds. At the beginning, we have Total-Cost resource budget for sampling. After we visit a vertex, its neighbors are all discovered and sampling strategy is used to decide which vertex is sampled next. Then this process is iterated until Total-Cost is used up in the sampling process. Different sampling methods differ in the size of seeds and sampling strategy. In this section, we describe several popular sampling methods in detail.

**Breadth-First Sampling (BFS):** BFS is a classical graph sampling algorithm, which has been widely applied in studying OSNs. For example, BFS is used for measurement and

topological characteristics analysis [2, 16] and user behavior analysis of OSNs [11, 23].

BFS starts from one random seed in the graph and aims at collecting the vertices which are close to the seed. Two queues are kept in the sampling process of BFS: queue *Sampled* stores sampled vertices, while queue *Waiting* stores vertices that will be sampled. Initially, the seed is put into queue *Waiting*. At each step, the first vertex  $v$  in queue *Waiting* is moved to queue *Sampled*, and all the neighbors of vertex  $v$  are added into queue *Waiting*, unless the neighbor is already in queue *Waiting* or queue *Sampled*. The process loops until the Total-Cost is used up. During the sampling process, once queue *Waiting* is empty, a random vertex will be selected and inserted into queue *Waiting*.

Recently, [9] points out that BFS is biased towards high-degree vertices. Moreover, due to our evaluation, BFS may only collect some local part of the graph and performs badly in disconnected or loosely connected graphs. Since BFS has the “ready to use” merit, it is still widely used for social graph sampling.

**Metropolis-Hasting Random Walk (MHRW):** Random Walk is another widely used graph sampling method. In RW, one random vertex is the initial seed and the sampling strategy is that we always choose one of the neighbors (vertex  $v$ ) of the current sampled vertex (vertex  $u$ ) as the next vertex, which spends Walk-Cost. Thus, the transition probability from  $u$  to  $v$  is

$$P_{u,v} = \begin{cases} 1/k_u & \text{if } v \text{ is } u\text{'s neighbor} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

RW is simple and practical but it is biased towards high-degree vertices. In RW, the next sampled vertex just depends on the current sampled vertex, thus the sampling process can be modeled as a Markov Chain. If the graph is undirected and non-bipartite, the Markov Chain is ergodic [14]. Therefore, the stationary probability of each edge in undirected graph is equal to  $1/|E|$  and each vertex is sampled with the probability  $\frac{k_u}{2|E|}$ . We can see that vertices with high-degrees are more likely to be sampled in RW.

To get rid of the bias towards high-degree vertices, the transition probabilities in RW should be modified appropriately and Metropolis-Hastings Random Walk is proposed in [9] to sample vertices uniformly. The Metropolis-Hastings algorithm is a Markov Chain Monte Carlo method to sample from a probability distribution whose direct sampling is difficult [15].

The algorithm of MHRW is as follows. Firstly, choose a vertex  $u$  as the initial seed. Secondly, select one neighbor (vertex  $v$ ) of vertex  $u$  and generate a random number  $p$  between 0 and 1 uniformly. If  $p < k_u/k_v$ ,  $v$  is chosen as the next sampled vertex, otherwise the walker still stays as vertex  $u$ . Then, the second step is iterated until Total-Cost is exhausted. The transition probability from  $u$  to  $v$  is

$$P_{u,v} = \begin{cases} \min(\frac{1}{k_u}, \frac{1}{k_v}) & \text{if } v \text{ is } u\text{'s neighbor} \\ 1 - \sum_{w \neq u} \min(\frac{1}{k_u}, \frac{1}{k_w}) & \text{if } v = u \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

From equation (2), we can see that the possibility to visit vertices with high degree is reduced and MHRW is proved to converge to uniform sampling, which means that each vertex in the whole graph is sampled with the probability  $1/|V|$  [9].

We should mention that one design assumption of MHRW is that the social graph is well connected. Besides, all duplicated vertices are valid and important in MHRW, since the duplication makes MHRW converge to uniform sampling, inherently.

**Random Jump (RJ):** Random Jump is a sampling method which is similar to RW. The difference between them is that the walker in RJ can jump to any random vertex in the graph with a fixed probability in each step. The benefit of jumping randomly is getting rid of the random walker being stuck in some locally well connected part of the graph. When the jump strategy is adopted, it will spend Jump-Cost. We should mention that RJ is still biased towards high-degree vertices, since if the jump strategy is not adopted, the vertices are selected just like RW.

Jump-Cost may be large in some OSNs. For instance, in Myspace, valid user-IDs are sparse in the space of user-IDs. Then, the size of sampled vertices obtained by RJ would be much smaller than the sample size obtained by RW. In practice, we usually prefer to get a large and fair sample, given the same Total-Cost. However, RJ may not perform well in respect of the size of sampled vertices.

**Frontier Sampling (FS):** All the sampling methods mentioned above are vertex sampling methods. Frontier Sampling [19], which performs multiple dependent random walkers in sampling graphs, is an edge sampling algorithm. Moreover, vertex sampling performs better than edge sampling in estimating small degree distribution [19], and in most OSNs vertices with small degree make up a major portion of the total graph. Thus we are focusing on improving the performance of vertex sampling methods in this paper. Nevertheless, we still introduce FS, which is a representative edge sampling method.

FS works in the following way. Firstly, multiple vertices are selected as the initial seeds. Secondly, choose a vertex  $u$  from the set of seeds with the probability proportional to its degree, i.e.  $P(u) \propto k_u$ . Thirdly, choose an edge  $(u, v)$  that starts from vertex  $u$  as a sample edge, and then replace  $u$  with  $v$  in the set of seeds. Repeat step 2 and step 3 until Total-Cost is exhausted.

FS is practical and samples edge uniformly, which indicates that the sample set itself is still biased towards high-degree vertices. Unlike MHRW, which aims at obtaining an unbiased sample directly, FS corrects the bias by specific estimators. In [19], estimators related to degree distribution and global clustering coefficient, which are based on sampled edges, are proposed. The results show that FS can achieve lower NMSE than RW for estimating degree distribution, especially in disconnected or loosely connected graphs.

RWRW is proposed in [9], and if it is treated as an edge sampling method, RWRW uses exactly the same estimators as FS. The comparisons between MHRW and RWRW [9] are also suitable for comparing MHRW and FS. Firstly, MHRW has the “ready to use” merit, since vertices are sampled uniformly. However, FS is biased towards high-degree vertices and requires re-weighting appropriately. Secondly, specific estimators should be built for estimating different graph properties. However, only estimators of degree distribution and global clustering coefficient are currently available and estimators for purely data-analytic procedures, such as hierarchical clustering or multidimensional scaling, is impossible to be constructed [9]. Thus, MHRW is more simple and versatile than FS in practice.

From the above description, we can conclude that BFS, RW and RJ do not converge to uniform sampling and are biased towards high-degree vertices. MHRW samples vertices uniformly, however, one design assumption of MHRW is that the social graph is well connected. FS, which is an edge sampling method, performs stably in disconnected or loosely connected graphs, but specific estimators should be built for estimating graph properties. Besides degree distribution and global clustering coefficient, estimators of any other properties are currently unavailable [19].

## 4. ALBATROSS SAMPLING

### 4.1 Workflow of Albatross Sampling

In this paper, we introduce random jump strategy into MHRW and propose an improved vertex sampling method named Albatross Sampling. The benefit of jump strategy in AS is to get rid of being stuck in some locally well connected part of the whole graph and gathers a “comprehensive” sample of the original graph.

To make this method simple, we fix the probability of jump in AS, i.e., if jump strategy is adopted, each vertex is sampled with the same probability  $1/|V|$  in this step. The algorithm of AS is described in Algorithm 1 and the code is available at [1]. Here,  $p$  represents the probability of jump and we choose  $p$  as 0.02 in our evaluation.

---

#### Algorithm 1 Albatross Sampling

---

```

cost  $\leftarrow$  0
sample set  $S \leftarrow$  empty
select a random vertex  $v$  as the initial seed
while cost < Total-Cost do
    generate  $\alpha$  from uniform distribution  $U[0, 1]$ 
    if  $\alpha < p$  then
        choose a new random vertex  $u$ 
        add  $u$  to the sample set  $S$ 
        if  $u$  is visited for the first time then
            cost  $\leftarrow$  cost+Jump-Cost
        else
            cost  $\leftarrow$  cost+Jump-Cost-1
        end if
         $v \leftarrow u$ 
    else
        select an edge  $(v, w)$  starting from  $v$  randomly
        generate  $\beta$  from uniform distribution  $U[0, 1]$ 
        if  $\beta < k_v/k_w$  then
            add  $w$  to the sample set  $S$ 
            if  $w$  is visited for the first time then
                cost  $\leftarrow$  cost+Walk-Cost
            end if
        else
            remain at  $v$ 
            add  $v$  to the sample set  $S$ 
        end if
    end if
end while

```

---

### 4.2 Unbiasness of Albatross Sampling

In AS, the next sampled vertex only depends on the current sampled vertex, thus the sampling process can be modeled as a Markov Chain and the transition probability from vertex  $u$  to vertex  $v$  is

$$P_{u,v} = \begin{cases} \min(\frac{1-p}{k_u}, \frac{1-p}{k_v}) + \frac{p}{|V|} & \text{if } v \text{ is } u\text{'s neighbor} \\ 1 - p - \sum_{w \neq u} \min(\frac{1-p}{k_u}, \frac{1-p}{k_w}) + \frac{p}{|V|} & \text{if } v = u \\ \frac{p}{|V|} & \text{otherwise} \end{cases}$$

In the above equation,  $p$  is the probability of jump in AS and  $|V|$  represents the total number of vertices in the original graph. In AS, the inherent Metropolis-Hasting algorithm guarantees that each vertex is sampled uniformly when jump strategy is not adopted. Then, when jump strategy is adopted, the probability of jumping to each vertex in the whole graph is equal. Therefore, the stationary probability of each vertex is equal to  $1/|V|$ .

Benefit from random jump strategy, AS avoids being trapped in locally well connected part of the social graph and converge to uniform sampling quickly. Due to our evaluation, AS estimates degree distribution with lower NMSE by consuming the same resource budget and converge quickly with smaller convergence time, even sampling disconnected or loosely connected graphs. Therefore, Albatross Sampling is a promising robust and effective vertex sampling method for social network analysis.

## 5. EVALUATION

In this section, we evaluate the performance of the sampling algorithms mentioned above. The comparison is focused on the performance of BFS, MHRW and AS. BFS is not an unbiased sampling method, nevertheless, we still compare with BFS, for its widely use in OSNs analysis [2, 11, 16, 23]. Here, we use two criterions to evaluate different sampling algorithms.

(1) **The accuracy of estimating degree distribution:** Normalized Mean Square Error (NMSE) is used for evaluating the estimating accuracy of these methods, which is also used in [19]. NMSE for degree  $k$  is defined as

$$NMSE(k) = \frac{\sqrt{E[(\hat{\theta}_k - \theta_k)^2]}}{\theta_k} \quad (3)$$

In this paper, we use NMSE as a criterion to evaluate the robustness of sampling methods. If a sampling algorithm achieves a lower NMSE in estimating degree distribution given the same sampling budget, it is regarded as more robust.

(2) **The converge rate of algorithms:** Firstly, we choose the evolution of estimating  $\theta_{10}$  for comparison, without loss of generality. Here,  $\theta_{10}$  means the fraction of vertices with in-degree (out-degree) less than or equal to 10, which is also used in [19]. Moreover, we give out the convergence time of different sampling methods. Similar to the mixing time of a Markov chain [13], we define *convergence time* as the cost  $c_0$  such that

$$|\hat{\theta}_k(c) - \theta_k| \leq 1/4 \quad (4)$$

for  $0 \leq k \leq \max(\text{in-degree}, \text{out-degree}), c_0 \leq c$

Convergence time reflects the least budget required for getting an acceptable estimation of the degree distribution.

### 5.1 Data Set

The data sets used for evaluation are Buzznet [3] and Berk-Stan [20], and their basic information is shown in Table 1. Buzznet is a photo, journal, and video-sharing social

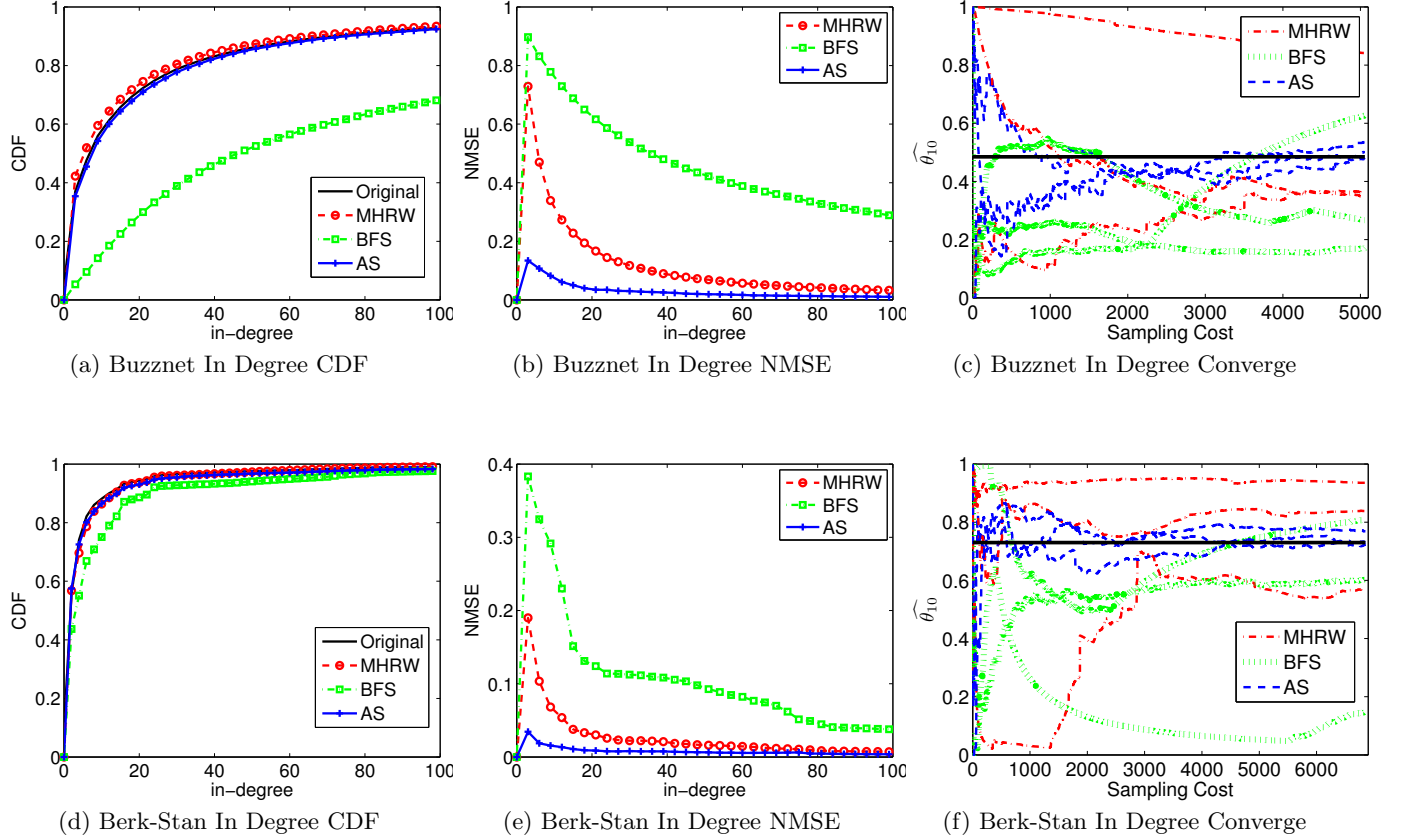


Figure 1: CDF, NMSE and Converge Curves for Buzznet(Cost= $|V|/20$ ) and Berk-Stan(Cost= $|V|/100$ )

media network where vertices represent users in the graph. In Buzznet, user A can subscribe to the updates of user B, which indicates a directed edge from A to B. Berk-Stan is the web graph of Berkeley and Stanford collected in 2002. Vertices represent pages from berkely.edu and stanford.edu domains and hyperlinks between them are represented as directed edges.

Table 1: Basic Information of Data Set

	Type	Nodes	Edges	SCC
Buzznet	Directed	101,169	4,284,534	0.944
Berk-Stan	Directed	685,230	7,600,595	0.489

In table 1, Strongly Connected Components (SCC) represents the fraction of number of vertices in the largest strongly connected component [5]. SCC shows the connectivity of a graph: if the value of SCC is smaller, the graph is more loosely connected. Thus, Buzznet is a tightly connected graph and Berk-Stan is a loosely connected graph. And we use both tightly and loosely connected graphs to compare these algorithms.

## 5.2 The Accuracy of Estimate

To make the analysis of degree distribution impressive, we plot the Cumulative Distribution Function (CDF) and Normalized Mean Square Error (NMSE) of degree distribution of the sampled vertices. In Figure 1 (a) and (b), the in-degree CDF and NMSE of Buzznet are presented. We

choose Total-Cost for sampling as 5% of the total number of vertices of Buzznet. In Figure 1 (d) and (e), we present the in-degree CDF and NMSE of Berk-Stan. We choose Total-Cost for sampling as 1% of the total number of vertices of Berk-Stan. Since vertices with in-degree less than 100 make up a major portion (more than 90%) of the total vertices in both these two datasets, comparisons of CDF and NMSE are focused on small degree vertices.

From these figures, we can see that BFS is biased towards high-degree vertices significantly, whose NMSE is much larger than that of MHRW and AS. Moreover, the CDF curves of MHRW and AS are both almost identical to the original CDF, but AS achieves smaller NMSE than MHRW. From the above description, we can conclude that AS is more robust for estimating degree distribution in both tightly and loosely connected graphs.

## 5.3 The Converge Rate of Algorithms

To evaluate comprehensively, we also compare the converge rate of these sampling methods. Figure 1 (c) and (f) show three sample paths of the evolution of  $\hat{\theta}_{10}$  as a function of sampling cost over Buzznet and Berk-Stan. And the black lines show the true value of  $\theta_{10}$ . We can see that two out of three sample paths in BFS underestimate  $\theta_{10}$ , which indicates that BFS is biased towards high-degree vertices. MHRW performs badly in these two graphs and none of three paths converge to the true value. However, all three paths of AS converge quickly and stably to the true value.

Moreover, in Table 2, we present the convergence time of different sampling methods in Buzznet and Berk-Stan. The convergence time of AS is 9.3% (Buzznet data set) and 5.1% (Berk-Stan data set) of the convergence time of BFS and 12.0% (Buzznet data set) and 7.1% (Berk-Stan data set) of the convergence time of MHRW. From these simulation results, we find that AS is more effective and reliable for converging quickly and stably.

**Table 2: Convergence Time**

	Buzznet	Berk-Stan
BFS	4371.2	6355.2
MHRW	3361.4	4652.5
AS	406.5	329.6

## 6. CONCLUSIONS AND FUTURE WORK

In this paper, we propose an improved vertex sampling algorithm named Albatross Sampling, which introduces random jump strategy into MHRW during the sampling process. Due to our evaluation, AS is more robust for estimating degree distribution with lower NMSE and more effective for converging more quickly with much smaller convergence time than MHRW and BFS, given the same sampling cost. Moreover, AS is reliable for sampling both tightly and loosely connected social graphs.

In the future, we will test on more large-scale social graphs and use AS to sample real OSNs. Also, we will evaluate the performance of these sampling algorithms in estimating other important graph properties, such as Assortativity [17] and Clustering Coefficient [21, 22]. Besides, one challenge for AS is that if valid user-IDs are quite sparse in the space of user-IDs, Jump-Cost would be too large. For instance, in December 2010, Facebook advertised over 500 million active users and each user is encoded by a 64-bit user-ID [10], resulting Jump-Cost is huge in Facebook. Proper improvement of AS should be considered for sampling social graphs with large Jump-Cost.

## 7. ACKNOWLEDGMENTS

This work is supported by the National Basic Research Program of China (No.2007CB310806) and the National Science Foundation of China (No.60850003, No.60473087).

## 8. REFERENCES

- [1] Project of Sampling Social Graphs. <http://code.google.com/p/sampling-social-graphs/>, 2011.
- [2] Y.-Y. Ahn, S. Han, H. Kwak, S. Moon, and H. Jeong. Analysis of Topological Characteristics of Huge Online Social Networking Services. In *Proc. of WWW*, 2007.
- [3] ASU. Social Computing Data Repository at Arizona State University. <http://socialcomputing.asu.edu/pages/home>.
- [4] F. Benevenuto, T. Rodrigues, M. Cha, and V. Almeida. Characterizing User Behavior in Online Social Networks. In *Proc. of ACM IMC*, 2009.
- [5] R. Diestel. *Graph Theory*. Springer-Verlag, 2010.
- [6] N. Eagle, A. S. Pentland, and D. Lazer. Inferring friendship network structure by using mobile phone data. *PNAS*, 106(36):15274–15278, August 2009.
- [7] Geek.com. Twitter reaches 200 million users and 110 million tweets per day. <http://www.geek.com/articles/news/twitter-reaches-200-million-users-and-110-million-tweets-per-day-20110120/>.
- [8] M. Gjoka. Measurement of online social networks. *UC Irvine PhD Thesis*, 2010.
- [9] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Walking in Facebook: A Case Study of Unbiased Sampling of Osns. In *Proc. of IEEE INFOCOM*, 2010.
- [10] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Practical recommendations on sampling osn users by crawling the social graph. *JSAC special issue on Measurement of Internet Topologies*, 2011.
- [11] H. Kwak, C. Lee, H. Park, and S. Moon. What is Twitter, a Social Network or a News Media? In *Proc. of WWW*, 2010.
- [12] J. Leskovec and C. Faloutsos. Sampling from Large Graphs. In *Proc. of ACM SIGKDD*, 2006.
- [13] D. A. Levin, Y. Peres, and E. L. Wilmer. *Markov Chains and Mixing Times*. American Mathematical Society, 2008.
- [14] L. Lovasz. Random walks on graphs: a survey. *Combinatorics*, 1993.
- [15] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equation of state calculation by fast computing machines. *J. Chem. Physics*, 106(21), August 1953.
- [16] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee. Measurement and Analysis of Online Social Networks. In *Proc. of ACM IMC*, volume 106, October 2007.
- [17] M. Newman. Assortative mixing in networks. *Phys. Rev. Lett.*, 89, 2002.
- [18] B. Ribeiro, W. Gauvin, B. Liu, and D. Towsley. On Myspace account spans and double Pareto-like distribution of friends. *UMass Amherst Technical Report*, (UM-CS-2010-001), January 2010.
- [19] B. Ribeiro and D. Towsley. Estimating and Sampling Graphs with Multidimensional Random Walks. In *Proc. of ACM IMC*, November 2010.
- [20] Stanford. Stanford large network dataset collection. <http://snap.stanford.edu/data/index.html/>.
- [21] T. Wang, Y. Chen, Z. Zhang, T. Xu, L. Jin, P. Hui, B. Deng, and X. Li. Understanding graph sampling algorithms for social network analysis. In *the 3rd ICDCS Workshop on Simplifying Complex Networks for Practitioners*, June 2011.
- [22] D. J. Watts and S. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393(6684), June 1998.
- [23] C. Wilson, B. Boe, A. Sala, K. P. N. Puttaswamy, and B. Y. Zhao. User Interactions in Social Networks and their Implications. In *Proc. of ACM EuroSys*, 2009.