

Quick detection of popular entities in large directed networks

ABSTRACT

In this paper, we address a problem of quick detection of popular entities in large online social networks. Practical importance of the problem is attested by a large number of companies that continuously collect and update statistics about popular entities. We suggest an efficient two-stage algorithm for solving this problem. For instance, our algorithm needs only one thousand API requests in order to find the top-50 most popular users in Twitter, a network with more than a billion of registered users. Our algorithm is easy to implement, it outperforms existing methods, and serves many different purposes, such as finding most popular users or most popular interest groups in social networks. An important contribution of this work is the analysis of the proposed algorithm using the Extreme Value Theory – a branch of probability that studies extreme events and properties of largest order statistics in random samples. Using this theory, we derive accurate predictions for the algorithm’s performance and show that the number of API requests for finding top- k most popular entities is sublinear in the number of entities. Moreover, we formally show that the high variability among the entities, expressed through heavy-tailed distributions, is the reason for the algorithm’s efficiency. We quantify this phenomenon in a rigorous mathematical way.

1. INTRODUCTION

In this paper, we propose a randomized algorithm for quick detection of popular entities in large online social networks. The entities can be, for example, users or interest groups, user categories, geographical locations, etc. For instance, one can be interested in finding out a list of Twitter users with many followers or Facebook interest groups with many members. Practical importance of the problem is attested by a large number of companies that continuously collect and update statistics about popular entities (*twittercounter.com*, *followerwonk.com*, *twitaholic.com*, *www.insidefacebook.com*, *yavkontakte.ru* just to name a few).

The problem at hand may seem trivial, if one assumes

that the network structure and the relation between entities are known. However, even then finding, for example, top- k in-degree nodes in a directed graph G of size N takes the time $O(N)$. For large networks, such linear complexity is already too high. In fact, for any practical purpose, it is much more valuable to find an approximate result in a sublinear time than an exact result in a linear time. Furthermore, the data of current social graphs is typically available only to the owners of the social network, and can be obtained by other interested parties only through API requests. The rate of allowed API requests is usually quite small. For instance, Twitter has the limit of one access per minute for one standard API account. Then, in order to crawl the entire network with more than 500 million users one need more than 950 years. Clearly, we would like to find most popular entities using only a small number of API requests.

Formally, the problem addressed in this paper is as follows. Let V be a set of entities, usually users, that can be accessed using API requests. Let W also be another set of entities (possibly equal to V). We represent V and W as vertices of a bipartite graph (V, W, E) , where a directed edge $(v, w) \in E$, with $v \in V$, $w \in W$, represents a relation between v and w . For instance, if $V = W$ is a set of Twitter users, then $(v, w) \in E$ may mean that v follows w , or that v retweeted a tweet from w . Note that any directed graph $G = (V, E)$ can be represented equivalently by the bi-partite graph (V, V, E) . One can also suppose that V is a set of users and W is a set of interest groups, while the edge (v, w) represents that user v belongs to group w . Our goal is to quickly find top- k in-degree entities in W . In this setting, throughout the paper, we use the terms ‘nodes’ and ‘entities’ interchangeably.

The algorithm proposed in this paper can detect popular entities with high precision using very small number of API requests. Most of our experiments are performed on the Twitter graph, because it is a good example of a huge network (billion of registered users) and very limited rate of requests to API. We use only 1000 API request to find top-50 Twitter users with very high precision. We also demonstrate the efficacy of our approach on the example of online social network (to be specified in the camera-ready version) which had, at the time of article preparation, more than 200 million registered users. We use our algorithm to quickly detect most popular interest groups in this social network. Experiments on random graph models show that our algorithm outperforms the baselines algorithms from [4] and [14]. Moreover, our algorithm can be used in a very general settings for finding most popular entities, while the baseline algorithms can only be use for finding nodes of largest de-

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.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

greens in directed ([14]) or undirected ([4]) graphs.

An important contribution of this work is the novel analysis of proposed algorithm using classical results of the Extreme Value Theory (EVT) – a branch of probability that studies extreme events and properties of largest order statistics in random samples. We refer to [8] for a comprehensive introduction to EVT. Specifically, we treat the largest in-degrees in W as high order statistics of a heavy-tailed distribution, and use EVT to obtain the limiting properties of these order statistics. This way we obtain statistical estimation of the magnitude of the largest in-degrees in W . Using these mathematical tools, we can, for instance, accurately predict the average fraction of correctly identified top-100 most followed users in Twitter using only the knowledge of top-20 degrees, which can be detected by our algorithm very quickly with practically 100% accuracy.

We derive the complexity of our algorithm in terms of the number of entities in W show that the complexity is sublinear if the in-degree distribution in W is heavy tailed. Intuitively, this should be the case because the high variability of the in-degrees implies that the largest entities have extremely large number of in-links and thus are easy to find. We formalize this argument using the EVT results.

The algorithm consists of two stages. The parameters of the algorithm, n_1 and n_2 , are the number of API requests used on the first and the second stage, respectively. The performance of the algorithm is very robust with respect of the parameters' values. We find the optimal scaling for n_1 and n_2 with respect to three measures of the algorithm performance: the average fraction of correctly identified top- k entities, the first-error index (the number of the highest statistics within top- k that was not included in the identified top- k list), and the the sum of incoming degrees of identified n_2 entities. Notice that for fixed n , the latter performance measure does not monotonically grows with n_2 because with small n_1 the number of links received from n_1 random users is a poor indication of the node's actual degree. This can be clearly seen in Figure 2 for the Twitter graph.

The rest of the paper is organized as follows. In Section 2, we give a short overview of the related work. In Section 3, we formally describe our algorithm. We empirically show the efficiency of our algorithm and compare it to baseline strategies in Section 4. We present a detailed analysis of the algorithm in Section 5 and evaluate its optimal parameters with respect to the above mentioned performance characteristics. Section 6 concludes the paper.

2. RELATED WORK

Over the last years data sets have become increasingly massive. For such large data any complexity higher than linear (in dataset size) is unacceptable, and even linear complexity may be too high. It is also well understood that an algorithm, which runs in sublinear time, cannot return an exact answer. In fact, such algorithms often use randomization, and then errors occur with positive probability. Nevertheless, in practice, a rough but quick answer is often more valuable than exact but computationally demanding solution. Therefore, sublinear time algorithms become increasingly important, and many studies of such algorithms have appeared in recent years (see, e.g., [10, 13, 15, 16]).

An essential assumption of this work is that the network structure is not available, and has to be discovered using the API requests. This setting is similar to on-line compu-

tations, when information is obtained and immediately processed while crawling the network graph (for instance the World Wide Web). There is a large body of literature where such on-line algorithms are developed and analyzed. Many of these algorithms are developed for computing and updating the PageRank vector [1, 6]. In particular, the algorithm recently proposed in [6] computes the PageRank vector in sublinear time. Furthermore, the probabilistic Monte Carlo methods [2, 11] allow to continuously update the PageRank as the structure of the Web changes.

Randomized algorithms are also used for discovering the structure of social networks. Often random walks are designed in such a way that the desired nodes are easily found. For example, in [12] an unbiased random walk, where each node is visited with equal probability, is constructed in order to find the degree distribution on Facebook. A different random walk is designed in [4] for finding nodes with largest degrees in undirected graphs. This random walk has jumps, so that it does not get stuck around just one hub, but unlike PageRank, its a stationary distribution completely defined by the nodes' degrees.

The problem of finding the most popular entities in large networks has been analyzed in several papers. In Section 4.3 we show that our algorithm outperforms two baselines: the random walk algorithm in [4], and the crawling algorithm in [14]. The latter algorithm [14] is designed to efficiently discover the correct set of pages with largest incoming degrees in a fixed network, and to track these pages over time when the network is changing. Their setting is different from ours in several aspects. For example, in our case we can use API to get indegree of any given item, while in the World Wide Web this information is not available. On the other hand, the algorithm in [14] is designed to discover the graph structure, and cannot be easily adopted for other tasks, such as finding most popular use categories or interest groups.

To the best our knowledge, this is the first work that presents and analyzes an efficient algorithm for retrieving most popular entities under realistic API constraints.

3. ALGORITHM DESCRIPTION

Recall that we consider a bipartite graph (V, W, E) , where V and W are sets of entities, and $(v, w) \in E$ represents a relation between the entities.

Let $n = n_1 + n_2$. Our algorithm has two stages, described below. See Algorithm 1 for the pseudocode.

First stage. We start by sampling uniformly at random a set A of n_1 entities (users, or nodes) $v_1, \dots, v_{n_1} \in V$. The nodes are sampled independently, so the same node may appear in A more than once, in which case we regard each copy of this node as a different node. Since multiplicities occur with very small probability this does not affect the efficiency of the algorithm but simplifies the implementation. For each node in A we record its out-neighbors in W . In practice, we bound the number of recorded out-links by the maximal number of id's that can be retrieved within one API request, thus the first stage uses exactly n_1 API requests.

Second stage. Let $S_w, w \in W$, be the number of nodes in A that have a (recorded) edge to w , and let w_i be the node in W with i -th largest values of S_w , so that $S_{w_1} \geq S_{w_2} \geq \dots \geq S_{w_N}$. Then we use another n_2 API requests to retrieve the actual in-degrees of the n_2 top-nodes w_1, \dots, w_{n_2} .

The set $\{w_1, w_2, \dots, w_{n_2}\}$ is supposed to contain nodes from W with large in-degrees. For example, if we are inter-

ested in top- k in-degree nodes in a directed graph, we hope to identify these nodes with high precision if k is significantly smaller than n_2 .

Algorithm 1: Find entities with large incoming degrees

input : Set of entities V of size M , set of entities W of size N , number of random nodes n_1 , number of candidate nodes n_2

output: Nodes $w_1, \dots, w_{n_2} \in W$, their degrees d_1, \dots, d_{n_2}

for $i \leftarrow 1$ **to** N **do**

$S[i] \leftarrow 0$;

for $i \leftarrow 1$ **to** n_1 **do**

$v \leftarrow \text{random}(M)$;

$F \leftarrow \text{OutNeighbors}(v)$;

foreach j **in** F **do**

$S[j] \leftarrow S[j] + 1$;

$w_1, \dots, w_{n_2} \leftarrow \text{Top}_{n_2}(S)$ // $S[w_1], \dots, S[w_{n_2}]$ are top n_2 values in S ;

for $i \leftarrow 1$ **to** n_2 **do**

$d_i \leftarrow \text{InDegree}(w_i)$;

4. EXPERIMENTS

4.1 Twitter graph

First, we show that our algorithm quickly finds the most popular users in Twitter graph. Formally, V is a set of Twitter users, $W = V$, and $(v, w) \in E$ iff v is a follower of w . Twitter is an example of a huge network with limited access to its structure. Information on the Twitter graph can be obtained via Twitter API. The standard rate of requests to API is one per minute. Every vertex has an id, which is an integer number starting from 12. The largest id of a user is $\sim 1460M$ (at the time when we performed the experiments). Due to such id assignment, a random user in Twitter can be easily chosen. The only problem is that some users in this range have been deleted, some are suspended, and therefore errors occur when addressing the id's of these pages. In our implementation we usually skip errors and assume that we do not spend resources on such nodes. The fraction of errors is $\approx 20\%$.

Given an id of a user, a request to API can return one of the following: i) the number of followers (indegree), ii) the number of followees (outdegree), or iii) at most 5000 id's of followers or followees. If a user has more than 5000 followees, then all their id's can be retrieved only by using several API requests. Instead, as described above, we record only the first 5000 of the followees and ignore the rest. This does not affect the performance of the algorithm because we record followees of randomly sampled users, and the fraction of Twitter users with more than 5000 followees, is small.

In order to obtain the ground truth, we first took $n_1 = n_2 = 500000$ and found top-1000 users with a very high precision. We used the obtained list for evaluating the performance of our algorithm.

Figure 1 shows the average fraction of correctly identified users from top- k for different k over 100 experiments, as a

function of n_2 , when $n = 1000$. Remarkably, we can find top-50 users with very high precision.

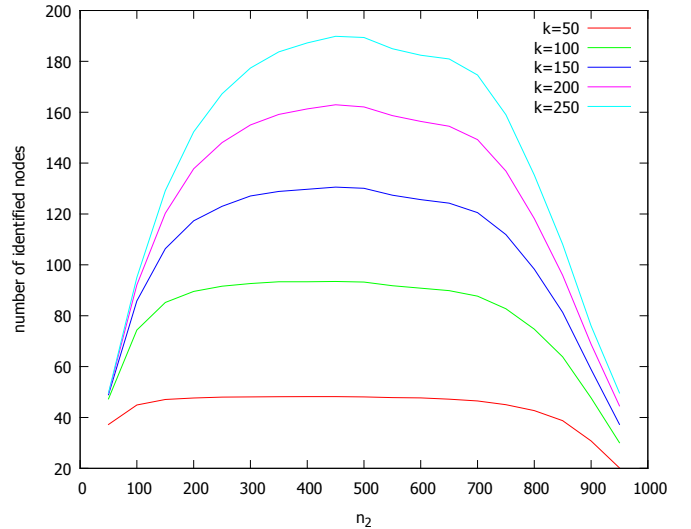


Figure 1: The number of correctly identified top- k most followed Twitter users as a function of n_2 , with $n = 1000$.

We have also looked at the first-error index – the position of the first mistake in the top- k list. Formally, if we correctly identified top- $(i-1)$ users, but did not find the i th user, then the first-error index is i . Again, we have averaged the results over 100 experiments. Results are shown in Figure 4 below (red line). Note that with only 1000 API requests we can correctly identify more than 50 users without any omission.

The sums of the degrees of the identified top- n_2 entities, with $n = 1000$, are depicted in Figure 2. Observe that here the optimal value of n_2 is larger than in two previously discussed metrics. Thus, in order to discover as many true in-links as possible, we may want to check more incoming degrees in the second stage of the algorithm, so that we have a large output list, but with less precision. We will discuss this in more detail in Section 5.3.

4.2 Finding largest interest groups

Let V be a set of users, W be a set of interest groups, and $(v, w) \in E$ iff v is a member of w .

We will demonstrate that our algorithm can find the most popular groups in a large social network with more than 200M registered users (to be specified in the camera-ready version). As for Twitter, information on the network under consideration can be obtained via API. Again, all users have ids: integer numbers starting from 1. Due to this id assignment, a random user in this network can be easily chosen. In addition, all interest groups also have their own id's.

We are interested in the following requests to API: i) given id of a user, return his or her interest groups, ii) given id of a group return its number of members. If a user decide to hide the list of groups, then an error occurs. The portion of such errors is $\approx 30\%$.

As before, first we used our algorithm with $n_1 = n_2 = 50000$ in order to find the most popular groups with high precision. Table 1 presents some statistics on the most popular groups. Then, we took $n_1 = 700$, $n_2 = 300$ and com-

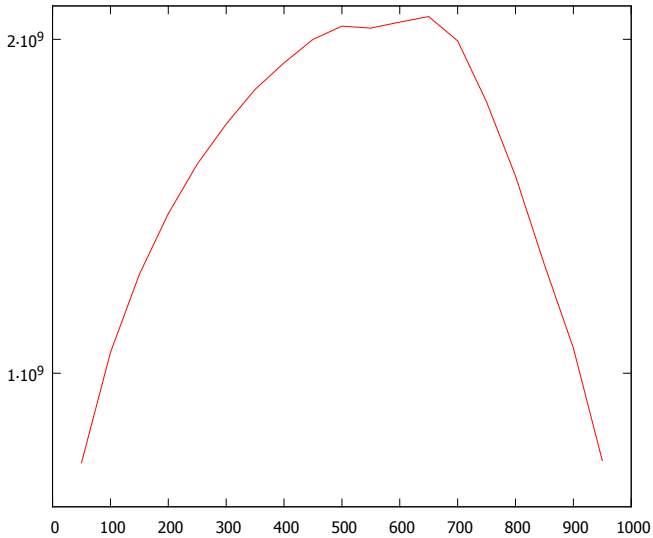


Figure 2: The sum of incoming degrees of identified users as a function of n_2 , $n = 1000$.

Table 1: The most popular groups

Rank	Number of participants	Topic
1	4,35M	humor
2	4,1M	humor
3	3,76M	movies
4	3,69M	humor
5	3,59M	humor
6	3,58M	facts
7	3,36M	cookery
8	3,31M	humor
9	3,14M	humor
10	3,14M	movies
100	1,65M	success

puted the fraction of correctly identified groups from top-100. Using only 1000 API requests, our algorithm identifies on average 73.2 from the top-100 interest groups (averaged over 25 experiments). The standard deviation is 4.6.

4.3 Comparison with baseline algorithms

In this section we compare our algorithm with the algorithms suggested in [4] and [14]. We start with the description of these algorithms.

Random walk based algorithm [4]. The algorithm in [4] is a randomized algorithm for undirected graphs that finds a top- k list of nodes with largest degrees in sublinear time. It is based on the random walk with uniform jumps, described by the following transition probabilities [5]:

$$p_{ij} = \begin{cases} \frac{\alpha/N+1}{d_i+\alpha}, & \text{if } i \text{ has a link to } j, \\ \frac{\alpha/N}{d_i+\alpha}, & \text{if } i \text{ does not have a link to } j, \end{cases} \quad (1)$$

where N is the number of nodes in the graph and d_i is the degree of node i . The parameter α controls how often the random walk makes an artificial jump. In [4] it is suggested to take the parameter α equal to the average degree

Algorithm 2: Random walk based algorithm

input : Graph G with N nodes, E edges, number of steps n , size of output list k , parameter α
output: Nodes v_1, \dots, v_k , their degrees d_1, \dots, d_k

```

 $v \leftarrow \text{random}(N)$ ;
 $F \leftarrow \text{Neighbors}(v)$ ;
 $D[v] \leftarrow \text{size}(F)$ ;
for  $i \leftarrow 2$  to  $n$  do
   $r \xleftarrow{\text{sample}} U[0, 1]$ ;
  if  $r < \frac{D[v]}{D[v]+\alpha}$  then
     $v \leftarrow \text{random from } F$ ;
  else
     $v \leftarrow \text{random}(N)$ ;
     $F \leftarrow \text{Neighbors}(v)$ ;
     $D[v] \leftarrow \text{size}(F)$ ;
 $v_1, \dots, v_k \leftarrow \text{Top}_k(D)$  //  $D[v_1], \dots, D[v_k]$  are top  $k$  values in  $D$ ;

```

in order to maximize the number of independent samples. Interestingly, this implies that the random walk, in stationarity, makes on average just one step between the jumps. With such choice of α the random walk method of [4] mimics most closely the suggested algorithm with independent sampling and exactly one step from entity in V to entity in W . We should note that the random walk method could be very valuable when the independent uniform sampling is expensive, for example, when the id space is very sparse.

The random walk keeps a candidate list of k nodes. Once a new node is discovered according to the transition probability (1), we check its degree and compare it with degrees of the nodes in the candidate list. If this newly discovered node has a degree larger than degrees of some nodes in the candidate list, the newly discovered node is inserted in the candidate list and a node with the smallest degree in the candidate list is pushed out. See Algorithm 2 for more detailed description. The algorithm can be run for a predefined number of steps or can be terminated according to one of the stopping criteria provided in [4].

Crawl-AI and Crawl-GAI [14]. At each step we consider one node and ask for its outgoing edges. At step n node j has its *apparent indegree* S_j , $j = 1, \dots, N$: the number of discovered edges pointing to this node. In Crawl-AI the next node to consider is a random node, with probability proportional to the apparent indegree. In Crawl-GAI, the next node is the node with the highest apparent indegree. After n steps we get a list of nodes with largest apparent indegrees. See Algorithm 3 for the pseudocode of the algorithm Crawl-GAI.

In the experiments of the present paper we take the same budget for all tested algorithms to compare their performance.

Note that we cannot compare the algorithms on the Twitter graph for several reasons. First, Algorithm 2 works only on undirected graphs. Second, in order to choose a random edge of a node, we need at least two request to API, to ask for followees and followers. Also, the random walk often hits nodes of high degree, and then many additional requests are needed to retrieve their followers and followees, because the

Algorithm 3: Crawl-GAI

input : Graph G with N nodes, number of steps n , size of output list k
output: Nodes v_1, \dots, v_k
for $i \leftarrow 1$ **to** N **do**
 $S[i] \leftarrow 0$;
for $i \leftarrow 1$ **to** N **do**
 $v \leftarrow \operatorname{argmax}(S[i])$;
 $F \leftarrow \operatorname{OutNeighbors}(v)$;
 foreach j **in** F **do**
 $S[j] \leftarrow S[j] + 1$;
 $v_1, \dots, v_k \leftarrow \operatorname{Top-k}(S)$ // $S[w_1], \dots, S[w_k]$ are top k values in S ;

Table 2: Number of correctly identified nodes from top-100 averaged over 100 experiments, $n = 1000$.

Algorithm	mean	standard deviation
Our (directed)	91.9	4.88
Crawl GAI (directed)	81.9	2.42
Crawl AI (directed)	82.9	2.38
Our (undirected)	97.9	1.71
Random walk (undirected)	60.7	4.76

number of id’s that can be obtained in one request is limited (5000 in Twitter). For example, we need 6K request to get the followers of a user with 30M followers. Algorithm 3 crawls only out-degrees, that are usually much smaller, but it can potentially suffer from the API constraints, for example, when in-degrees and out-degrees are dependent.

Therefore, in order to compare Algorithms 1–3, we have generated a random directed graph according to the configuration model (see [7]). Our artificial graph has 1M nodes, 6M edges, and the parameter of the power law degree distribution is 2. This directed graph is used to compare our algorithm to Crawl-AI and Crawl-GAI. In order to compare our method to the random walk based algorithm, we treat the generated graph as undirected. As prescribed by [4], we took α slightly smaller than the average degree in the graph (in our case $\alpha = 10$) and we considered a random walk with 1000 steps.

For the algorithm suggested in this paper we took $n_1 = 700$, $n_2 = 300$. The results of comparison can be seen in Table 2.

We expect our algorithm with $n_1 = 1000$ to be close to Crawl-GAI. Indeed, in the directed case our algorithm with $n_1 = 1000$ identifies 81.4 nodes from top-100 on average (this number is not presented in the table). Further improvement of our algorithm over the baselines is obtained because of the right balance between n_1 and n_2 .

5. ANALYSIS OF THE ALGORITHM

In this section, we present the theoretical analysis of Algorithm 1. The goal of this analysis is: 1) to mathematically justify our suggested two-steps procedure; 2) to prove that the total number of API requests, n , scales sublinearly with the network size, N ; 3) to find the optimal scaling of n_1 and n_2 (the number of API requests in the first and the second

stage of the algorithm) with respect to n .

We number the nodes in W by $1, 2, \dots, N$ according to the number of incoming links, from most popular to least popular. As prescribed by Algorithm 1, we pick n_1 nodes in V uniformly at random. The first important observation is that S_j follows a binomial distribution. Indeed, let F_j be the unknown random in-degree of node $j \in W$, so that $F_1 \geq F_2 \geq \dots \geq F_N$. Then, if we label all nodes from V that have a edge to j (we call such nodes followers of j), then S_j is exactly the number of labeled nodes in a random sample of n_1 nodes, so its distribution is $\operatorname{Binomial}\left(n_1, \frac{F_j}{N}\right)$. Hence, we have

$$\mathbb{E}(S_j|F_j) = n_1 \frac{F_j}{N}, \quad \operatorname{Var}(S_j) = n_1 \frac{F_j}{N} \left(1 - \frac{F_j}{N}\right). \quad (2)$$

For the top nodes with large F_j this distribution can be approximated with the Poisson distribution $\operatorname{Poisson}\left(\frac{n_1 F_j}{N}\right)$.

5.1 Candidate list

The quality of the top- k lists produced by Algorithm 1 is defined by the events whether or not the value of S_j , $j = 1, \dots, k$, is among the top- n_2 values of S_1, S_2, \dots, S_N , obtained in the first stage of the algorithm. This is justified by the intuition that if $F_j > F_l$, then we are likely to see $S_j > S_l$. Note, however, that the case when S_j is as small as 1, the event $1 = S_j > S_l = 0$ is not informative.

EXAMPLE 1. *Let us take $n_1 = n_2 = 500$ in the case of the Twitter graph. Then the average number of nodes i among the top-10000 with $S_i = 1$ is already*

$$\sum_{i=1}^{10^4} P(S_i = 1) \approx \sum_{i=1}^{10^4} \frac{500 F_i}{5 \cdot 10^8} e^{-500 F_i / 5 \cdot 10^8} = 2539.1,$$

hence, many more than n_2 nodes will have $S_i = 1$ and can make it to the top n_2 values of S_1, S_2, \dots, S_N only with a small probability.

Motivated by the above considerations, we formulate our approach in terms of a statistical test as follows. Let our data be S_1, S_2, \dots, S_N . We assume that the observations are realizations of independent Poisson random variables with parameters $n_1 F_1 / N, n_1 F_2 / N, \dots, n_1 F_N / N$. For the two numbers $j, l \in 1, \dots, N$, we test the null-hypothesis $H_0 : F_j \leq F_l$ against the alternative $H_1 : F_j > F_l$. Let $S_{i_1} \geq S_{i_2} \geq \dots \geq S_{i_{n_2}}$ be the top- n_2 order statistics of S_1, \dots, S_N obtained by Algorithm 1. Then the first stage of the algorithm is equivalent to rejecting $H_0 : F_{i_j} \leq F_{i_{n_2}}$ for $j = 1, \dots, n_2 - 1$ such that

$$S_{i_j} > \max\{S_{i_{n_2}}, 1\}. \quad (3)$$

Here the strict inequality is necessary to guarantee that i_j is on the top- n_2 list after the first stage of the algorithm. If H_0 is rejected, then the actual degree of entity i_j will be retrieved in the second stage of the algorithm.

Note that in contrast to the classical hypothesis testing, here we do not draw the conclusions solely from the observed random data S_1, S_2, \dots, S_N but we obtain the true values of the parameters in the second stage of the algorithm. Hence, if we use $S_{i_{n_2}}$ as a proxy for S_{n_2} , then, given F_1, F_2, \dots, F_N , the quality of the top- k list is expressed as the power of

the test as follows:

$$\begin{aligned}
& P(\text{node } j \text{ is found} | F_j, F_{n_2}) \\
&= P(S_j > \max\{S_{i_{n_2}}, 1\} | F_j, F_{i_{n_2}}) \\
&\approx P(S_j > \max\{S_{n_2}, 1\} | F_j, F_{n_2}) \\
&\approx \sum_{s=0}^{\infty} e^{-\frac{n_1 F_{n_2}}{N}} \frac{(n_1 F_{n_2})^s}{N^s s!} \sum_{r > \max\{s, 1\}} e^{-\frac{n_1 F_j}{N}} \frac{(n_1 F_j)^r}{N^r r!} \\
&=: P_j(n_1), \quad j = 1, \dots, k.
\end{aligned} \tag{4}$$

5.2 Performance criteria

The main constraint of Algorithm 1 is the number of API requests that we can use. In order to measure the performance of the algorithm, we propose three objectives, described formally in this section.

The first objective is the average number of correctly identified top- k nodes. This is defined in the same way as in [3]:

$$\begin{aligned}
& E[\text{fraction of correctly identified top-}k \text{ entities}] \\
&= \frac{1}{k} \sum_{j=1}^k P(\text{node } j \text{ is found} | F_j, F_{n_2}) \approx \frac{1}{k} \sum_{j=1}^k P_j(n_1). \tag{6}
\end{aligned}$$

The second objective is the first-error index, which is equal to i if the top $(i-1)$ entities are identified correctly, but the top- i entity is not identified. If all top- n_2 entities are identified correctly, we set the first-error index equal to $n+1$. Using that for a discrete random variable X with values $1, 2, \dots, k+1$ holds $E(X) = \sum_{l=0}^k P(X > l)$, we obtain the average first-error index as follows:

$$\begin{aligned}
& E[\text{1st-error index}] = \sum_{l=0}^{n_2} P(\text{1st-error index} > l) \\
&= \sum_{j=1}^{n_2+1} \prod_{l=1}^{j-1} P(S_j > \max\{S_{i_{n_2}}, 1\} | F_j, \dots, F_{i_{n_2}}) \\
&\approx \sum_{j=1}^{n_2} \prod_{l=1}^{j-1} P_l(n_1). \tag{7}
\end{aligned}$$

Finally, our last objective is the sum of the identified top- n_2 degrees, that can be written in a very simple form:

$$U := [\text{sum of identified } n_2 \text{ degrees}] = \sum_{l=1}^{n_2} F_{i_l}. \tag{8}$$

5.3 EVT performance predictions

In order to compute the values in (6), (7), we need to make assumptions on the top- n_2 in-degrees of entities in W : F_1, F_2, \dots, F_{n_2} . To this end, we employ the quantile estimation techniques from the Extreme Value Theory (EVT).

In most social networks the degrees of the entities show a great variability. This is often modeled using power laws, although it has been often argued that classical Pareto distribution does not always fit the observed data. In our analysis we assume that the incoming degrees of the entities in W are independent random variables following a *regularly varying* distribution G :

$$1 - G(x) = L(x)x^{-1/\gamma}, \quad x > 0, \tag{9}$$

where $L(\cdot)$ is a slowly varying function, that is,

$$\lim_{x \rightarrow \infty} L(tx)/L(x) = 1, \quad t > 0$$

($L(\cdot)$ can be, for example, a constant or a logarithm). We note that (9) describes a broad class of heavy-tailed distributions, for which the EVT arguments presented below are valid, without imposing the rigid Pareto assumption.

Observe that F_1, F_2, \dots, F_N are the order statistics of G . Assume now that we know the correct values of the top- m nodes, $m < k$. This is plausible because, for instance, in Twitter, with $n = 1000$, the top-50 nodes are identified with a very high precision, see Figure 1. Then, in order to estimate the value of γ , we can use the classical Hill's estimator $\hat{\gamma}$, based on the top- m order statistics:

$$\hat{\gamma} = \frac{1}{m-1} \sum_{i=1}^{m-1} \log(F_i) - \log(F_m). \tag{10}$$

Next, we use the quantile estimator, given by formula (4.3) in [9], but we replace their two-moment estimator by the Hill's estimator in (10). This is possible because both estimators are consistent (under slightly different conditions). Under the assumption $\gamma > 0$, we have the following estimator f_j for the $(j-1)/N$ -th quantile of G :

$$f_j = F_m \left(\frac{m}{j-1} \right)^{\hat{\gamma}}, \quad j > 1, j \ll N. \tag{11}$$

We propose to use f_j as a prediction of F_j .

Note that our argument is inspired but not entirely justified by [9] because the consistency of the proposed quantile estimator (11) is only proved for $j < m$, while we want to use it for $j > m$. However, in the experiments we observe that expressions (6) and (7) are very robust with respect to the estimated values F_1, \dots, F_{n_2} . Moreover, $\hat{\gamma}$ increases with m , and it is easy to see that with smaller $\hat{\gamma}$ the predictions of the algorithm performance are more conservative.

In Figure 3 we compare the true fraction of the correctly identified top- k followed Twitter users to the performance prediction (6) for $n = 1000$ and $k = 100$. The magenta line shows the prediction for the fraction of correctly identified nodes in (6), where we used the correct values of F_1, F_2, \dots, F_{n_2} . The green line represents the results for the estimated values of F_1, \dots, F_k and F_{n_2} , based on the true values of the top-20 degrees. We see that it is very close to the magenta line, which is based on the true values of the degrees.

Similarly, we use formula (7) and the estimator (11) in order to provide the prediction of the first-error index. The results are given in Figure 4. We see again that the EVT predictions are more pessimistic than the experimental results, so we find the lower bound for the algorithm's actual performance. Note also that the shape of the plot and the optimal value of n_2 have been captured correctly by both predictors.

It is also clear that there is a principal difficulty in finding similar analytical predictions for the objective U in (8) because it is based not on the *actual* degrees F_1, F_2, \dots , but on the degrees $F_{i_1}, F_{i_2}, \dots, F_{i_{n_2}}$, where $S_{i_1} \geq S_{i_2} \geq \dots \geq S_{i_N}$ are the order statistics of the S_j 's. The exact expressions for such order statistics are rather messy. However, we can get some insight in the behavior of U in Figure 2. Indeed, clearly, the sum of correct top- n_2 degrees, $\sum_{i=1}^{n_2} F_{i_i}$, is an increasing function of n_2 . Moreover, if we use the estimator (11), then we observe that the largest values of F_j 's

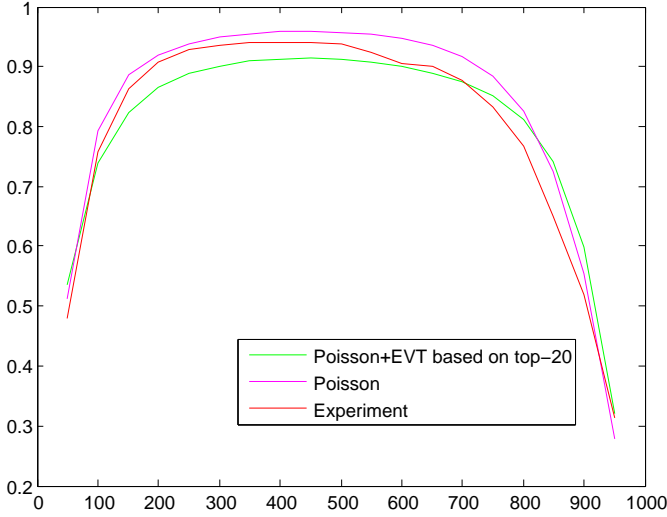


Figure 3: Fraction of correctly predicted nodes out of top-100 as a function of n_2 , with $n = 1000$: experiments (red); prediction (5) based on the true values of the degrees (magenta); prediction (5) based on top- m degrees and estimator (11) with $m = 20$, $\hat{\gamma} = 2.2$ (green).

are of the same order of magnitude:

$$F_j/F_i \approx \left(\frac{l-1}{j-1} \right)^{\hat{\gamma}}.$$

Thus, as long as n_1 large enough so that a large entity j receives large S_j , we have that U is comparable to $\sum_{i=1}^{n_2} F_i$, and hence U increases in n_2 . However, as n_1 becomes smaller, then small entities will constitute a large proportion of the set $\{i_1, i_2, \dots, i_{n_2}\}$. For example, if $n_2 = 800$, $n_1 = 200$, then we obtain, for the true values of in-degrees in Twitter graph with $N \approx 500M$:

$$\sum_{i=1}^{800} P(S_i > 1) \approx 280.9,$$

thus on average about 520 out of the top-800 nodes will be undistinguishable from other, much smaller nodes (see Example 1). Moreover, in this case

$$\sum_{i=1}^{10^5} P(S_i > 1) \approx 485.18,$$

thus, on average, more than 300 nodes will be included into $\{i_1, \dots, i_{800}\}$ essentially on a random basis. Since large majority of the nodes has very small degrees, this will drastically affect the magnitude of U . This is exactly what we observe in Figure 2.

5.4 Optimal scaling for n_1 and n_2

In this section our goal is to find the ratio n_2 to n_1 which maximizes the performance of the Algorithm 1. For simplicity, as a performance criterion we consider the fraction of correctly identified nodes from top- k in (6):

$$\frac{1}{k} \sum_{j=1}^k P_j(n_1) \rightarrow \max.$$

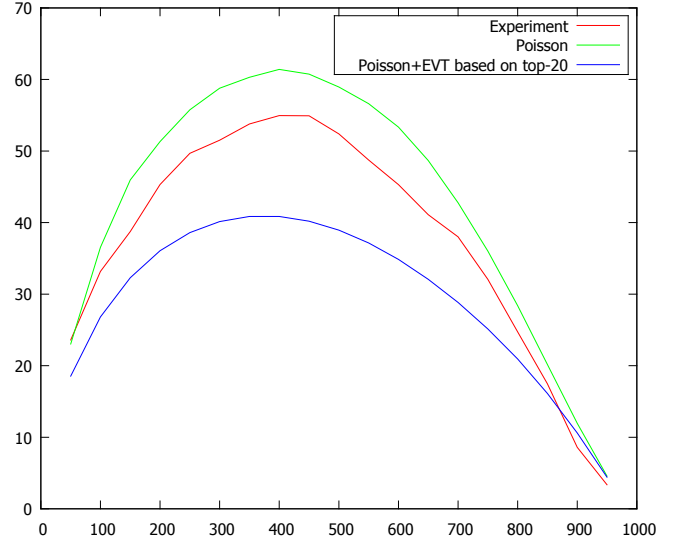


Figure 4: The position of first error as a function of n_2 , with $n = 1000$.

We start with analyzing the optimal scaling for n_1 . Intuitively, after the first stage of the algorithm, only $O(n_1)$ nodes j will have $S_j > 1$, and thus there is no need to check more than $n_2 = O(n_1)$ nodes in the second stage, which implies that n_1 should grow at least proportionally to n . This is formalized in the next proposition.

PROPOSITION 1. *It is optimal to choose $n = O(n_1)$.*

PROOF. Let J be a randomly chosen node, and J_l , $l = 1, \dots, n_1$ be independent realizations of J in the first stage of Algorithm 1. Denote by M the maximal number of neighbors that a given API allows to retrieve. The first stage of the algorithm returns a list of candidate nodes, for which we require $S_j > 1$. Observe that the number of such nodes is bounded by

$$U := \frac{1}{2} \sum_{l=1}^{n_1} \max\{M, \text{out-degree}(J_l)\}.$$

Assuming that the out-degrees of each node are independent, we obtain that

$$E(U) = \frac{1}{2} n_1 E(\max\{M, \text{out-degree}(J)\}),$$

$$\text{Var}(U) = \frac{1}{4} n_1 \text{Var}(\max\{M, \text{out-degree}(J)\}).$$

Note that the API restriction simplifies the derivation because the variance of $\max\{M, \text{out-degree}(J)\}$ is finite. The formal argument for $M = \infty$ and infinite variance of out-degrees will be similar but requires some more work. Using, e.g. Chernoff bound or Chebyshev bound we obtain that $P(U > E(U)(1 + \varepsilon)) \rightarrow 0$ as $n_1 \rightarrow \infty$. Thus, the number of nodes j with $S_j > 1$ is at most $O(n_1)$ with high probability, so we choose $n_2 = O(n_1)$ which results in $n = O(n_1)$. \square

Note that if n is large enough, then the top nodes (first, second, etc.) can be found with very high probability. Figure 1 shows that if $n = 1000$, then for a wide range of n_2 the fraction of correctly identified nodes from top-50 is the

same. As k grows, the optimization becomes much more important. Motivated by this observation, we maximize the value $P_k(n_1)$. We prove the following theorem.

THEOREM 1. *Assume that $k = o(n)$ as $n \rightarrow \infty$. The maximizer n_2^* of probability $P_k(n - n_2)$ is close to the maximal root of the equation*

$$\frac{1}{3\gamma k^\gamma} x^{\gamma+1} + x - n = 0, \quad (12)$$

that is,

$$n_2^* = x(1 + o(1)), \quad \text{as } k/n_2^* \rightarrow 0.$$

If in addition $n_2^* = o(n)$ as $n \rightarrow \infty$, then n_2^* can be given in a closed-form asymptotic expression

$$n_2 = (3\gamma k^\gamma n)^{\frac{1}{\gamma+1}} + o(n^{\frac{1}{\gamma+1}}).$$

PROOF. Consider first an extreme regime: $x = O(k)$. Thus, we exclude the regime $n - x = o(n)$. Consequently, $n_1 \rightarrow \infty$ as $n \rightarrow \infty$ and we can apply the following normal approximation

$$\begin{aligned} P_k(n_1) &\approx P\left(N\left(\frac{n_1(F_k - F_{n_2})}{N}, \frac{n_1(F_k + F_{n_2})}{N}\right) > 0\right) \\ &= P\left(N(0, 1) > -\sqrt{\frac{n_1}{N}} \frac{F_k - F_{n_2}}{\sqrt{F_k + F_{n_2}}}\right). \end{aligned} \quad (13)$$

(A completely formal justification can be given by the Berry-Esseen theorem.) Thus, in order to maximize the above probability, we need to maximize $\sqrt{\frac{n_1}{N}} \frac{F_k - F_{n_2}}{\sqrt{F_k + F_{n_2}}}$. From EVT it follows that F_k decays as $k^{-\gamma}$. So, we can maximize

$$\frac{\sqrt{n_1}(k^{-\gamma} - n_2^{-\gamma})}{\sqrt{k^{-\gamma} + n_2^{-\gamma}}}. \quad (14)$$

Now if $x = O(k)$, $\sqrt{n-x} = \sqrt{n}(1 + o(1))$, and the maximization of (14) mainly depends on the remaining term in the product, which is an increasing function of n_2 . This suggests that n_2 has to be chosen considerably greater than k . Hence, we proceed assuming the only interesting asymptotic regime where $k = o(n_2)$. In this asymptotic regime, we can simplify (14) as follows:

$$\begin{aligned} &\frac{\sqrt{n-x}(k^{-\gamma} - x^{-\gamma})}{\sqrt{k^{-\gamma} + x^{-\gamma}}} = \\ &\frac{1}{k^{\gamma/2}} \sqrt{n-x} \left(1 - \frac{3}{2} \left(\frac{k}{x}\right)^\gamma\right) + o\left(\left(\frac{k}{x}\right)^\gamma\right). \end{aligned}$$

Next, we differentiate the function

$$f(x) := \sqrt{n-x} \left(1 - \frac{3}{2} \left(\frac{k}{x}\right)^\gamma\right)$$

and set the derivative to zero. This results in equation (12). If we assume further that $n_2^* = o(n)$, then only the highest order term will remain in (12) and we immediately obtain the following approximation

$$n_2 = (3\gamma k^\gamma n)^{\frac{1}{\gamma+1}} + o(n^{\frac{1}{\gamma+1}}).$$

□

For example, for $n = 1000$, $k = 100$, and $\gamma = 0.35$ we get $n_2 \approx 570$.

5.5 Sublinear complexity

The normal approximation (13) immediately implies the following proposition.

PROPOSITION 2. *For large enough n_1 , the inequality*

$$\sqrt{\frac{n_1}{N}} \frac{F_k - F_{n_2}}{\sqrt{F_k + F_{n_2}}} \geq x_{1-\varepsilon}$$

guarantees that on average we can find the fraction $1 - \varepsilon$ of top- k nodes in W .

For the inequality in (2) to hold, it is necessary that $\sqrt{n_1}(F_k - F_{n_2})$ is at least of the same order of magnitude as $N\sqrt{F_k + F_{n_2}}$. Moreover, it follows from Proposition 1 that $n = O(n_1)$, and thus the complexity n of the algorithm is defined by n_1 . In the theorem below we use the results from Extreme Value Theory to show that n_1 scales sublinearly with N .

Theorem 1, and estimator (11), we can already provide a rough indication of the number of API request we need to use. Indeed, $k > m$, rough estimation with $n - n_2 \approx n$ and $F_k \gg F_{n_2}$ gives

$$n \geq \frac{Nx_{1-\varepsilon}^2 k^{\hat{\gamma}}}{F_m m^\gamma}. \quad (15)$$

For finding top-100 most followed users on Twitter with good precision, this will result in about 5000 of API requests (with $N = 500M$, $m = 20$, $k = 100$, $x_{1-\varepsilon} \approx 2$, $\hat{\gamma} = 2.2$).

For a better result, we may take into account the value of n_2 , and substitute the value $n_2 = (3k^\gamma n \hat{\gamma})^{\frac{1}{\gamma+1}}$ obtained in Proposition 2:

$$\begin{aligned} &\frac{k^{-\hat{\gamma}/2}}{2\sqrt{n}} \left(2n - (3k^\gamma n \hat{\gamma})^{\frac{1}{\gamma+1}}\right) \left(1 - \frac{3}{2} (3k^\gamma n \hat{\gamma})^{\frac{-\hat{\gamma}}{\gamma+1}} k^{\hat{\gamma}}\right) \\ &\geq x_{1-\varepsilon} \sqrt{\frac{N}{F_m m^\gamma}}. \end{aligned}$$

From (15) we can also already anticipate that n is sublinear in N because $F_m m^\gamma$ grows with N . This argument is formalized in Theorem 2 below.

Notice that, interestingly, the obtained complexity is in terms of the cardinality of W , not V . In particular, this makes the problem of finding popular groups easier than the problem of finding popular users.

THEOREM 2. *If the in-degrees of the nodes are independent realizations of a regularly varying distribution G with exponent $1/\gamma$ as defined in (9), and $F_1 \geq F_2 \geq \dots \geq F_N$ are their order statistics. Let $(a_N)_{N \geq 1}$, $(b_N)_{N \geq 1}$ be sequences such that*

$$\lim_{N \rightarrow \infty} N(1 - G(a_N x + b_N)) = (1 + \gamma x)^{-1/\gamma}.$$

Then Algorithm 1 finds $(1 - \varepsilon)$ of the top- k nodes with high probability in

$$n_1 = O(N/a_N),$$

of API requests. In particular, n scales sublinearly in N , and

$$\log(n_1) = (1 - \gamma) \log(N).$$

PROOF. For a regularly varying G , Theorem 2.1.1 in [8] can be applied, and thus for any finite m

$$\left(\frac{F_1 - b_N}{a_N}, \dots, \frac{F_m - b_N}{a_N}\right)$$

converges in distribution, as $N \rightarrow \infty$, to

$$\left(\frac{E_1^{-\gamma} - 1}{\gamma}, \dots, \frac{(E_1 + \dots + E_m)^{-\gamma} - 1}{\gamma} \right),$$

where E_i 's are independent exponential random variables with parameter 1. This implies, in particular, that $a_N/b_N = O(1)$ and that for large enough N and any $\varepsilon > 0$, there exist l_i, u_i such that $P[l_i a_N \leq F_i \leq u_i a_N] > 1 - \varepsilon$. It follows that for fixed k

$$\sqrt{\frac{n_1}{N}} \sqrt{F_k} = O(1)$$

with high probability when $n_1 = O(N/a_N)$, and the first statement of the theorem follows because $k = o(n_2)$ implying that $F_{n_2} = o(F_k)$. In particular, if G is a Pareto distribution, $1 - G(x) = Cx^{-1/\gamma}$, $x \geq x_0$, then

$$a_N = \gamma C^\gamma N^\gamma, \quad b_N = C^\gamma n^\gamma.$$

For a general regularly varying distribution in (9) the slowly varying function will influence a_N but the logarithmic asymptotics of a_N will be still determined by the power law:

$$\log(a_N) = \gamma \log(N),$$

which gives the result. \square

6. CONCLUSION

We proposed a randomized algorithm for quick detection of popular entities in large online social networks whose architecture has underlying directed graphs. Examples of social network entities are users and interest groups. We have analyzed the algorithm with respect to three criteria and compared with two baseline methods. Our analysis demonstrates that the algorithm has nonlinear complexity on networks with heavy-tailed in-degree distribution and that the performance of the algorithm is robust with respect to the values of its few parameters. The algorithm outperforms the two baseline methods and has much wider applicability. An important ingredient of our analysis is substantial use of the extreme value theory. The extreme value theory is not so well known in computer science and sociology but appears to be a very useful tool in the analysis of social networks. We feel that our work could be a good reference point for other researchers to start applying EVT in social network analysis. We have validated our theoretical results on two very large online social networks.

We see several extensions of the present work. A top list of popular entities is just one type of properties of social networks. We expect that our approach based on extreme value theory and using referral links can be extended to infer and to analyze other properties such as power law index and the tail, network functions and network motifs, degree-degree correlation. It will be very interesting and useful to develop quick and effective statistical tests to check for network assortativity and presence of heavy tails.

Since our approach requires very small numbers of API accesses, we believe that it will trace well network changes. Of course, a formal justification of the algorithm applicability for dynamic networks is needed.

7. REFERENCES

[1] S. Abiteboul, M. Preda, and G. Cobena. Adaptive on-line page importance computation. *Proceedings of*

the 12-th International World Wide Web Conference, 2003.

[2] K. Avrachenkov, N. Litvak, D. Nemirovsky, and N. Osipova. Monte carlo methods in pagerank computation: When one iteration is sufficient. *SIAM J. Numer. Anal.*, 45(2):890–904, 2007.

[3] K. Avrachenkov, N. Litvak, D. Nemirovsky, E. Smirnova, and M. Sokol. Quick detection of top-k personalized pagerank lists. In *Proceeding on the 8th Workshop on Algorithms and Models for the Web Graph, WAW 2011*, pages 50–61. Springer, 2011.

[4] K. Avrachenkov, N. Litvak, M. Sokol, and D. Towsley. Quick detection of nodes with large degrees. In *Proceeding on the 9th Workshop on Algorithms and Models for the Web Graph*, pages 54–65. Springer, 2012.

[5] K. Avrachenkov, B. Ribeiro, and D. Towsley. Improving random walk estimation accuracy with uniform restarts. In *Proceeding on the 7th Workshop on Algorithms and Models for the Web Graph, WAW 2010*, pages 98–109. Springer, 2010.

[6] C. Borgs, M. Brautbar, J. Chayes, and S.-H. Teng. A sublinear time algorithm for pagerank computations. *Lecture Notes in Computer Science*, 7323:41–53, 2012.

[7] T. Britton, M. Deijfen, and A. Martin-Löf. Generating simple random graphs with prescribed degree distribution. *J. Stat. Phys.*, 124(6):1377–1397, 2006.

[8] L. De Haan and A. Ferreira. *Extreme value theory*. Springer, 2006.

[9] A. L. M. Dekkers, J. H. J. Einmahl, and L. de Haan. A moment estimator for the index of an extreme-value distribution. *The Annals of Statistics*, 17(4):1833–1855, 1989.

[10] E. Fischer. The art of uninformed decisions: A primer to property testing. *Bull. Eur. Assoc. Theor. Comput. Sci. EATCS*, 75:97–126, 2001.

[11] D. Fogaras, B. Rácz, K. Csalogány, and T. Sarlósa. Towards scaling fully personalized pagerank: Algorithms, lower bounds, and experiments. *Internet Mathematics*, 2(3):333–358, 2005.

[12] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Walking in facebook: A case study of unbiased sampling of osns. *Proceedings of IEEE INFOCOM'10*, 2010.

[13] O. Goldreich. Combinatorial property testing: A survey. *Randomization Methods in Algorithm Design, DIMACS Ser. Discrete Math. Theoret. Comput. Sci.*, pages 45–60, 1998.

[14] R. Kumar, K. Lang, C. Marlow, and A. Tomkins. Efficient discovery of authoritative resources. *IEEE 24th International Conference on Data Engineering*, pages 1495–1497, 2008.

[15] R. Rubinfeld and A. Shapira. Sublinear time algorithms. *SIAM J. Discrete Math.*, 25(4):1562–1588, 2011.

[16] M. Sudan. Invariance in property testing. *Property Testing: Current Research and Surveys, O. Goldreich, ed., Lecture Notes in Comput. Sci.*, pages 211–227, 2010.