

# The Very Small World of the Well-connected

Xiaolin Shi  
Department of EECS  
University of Michigan  
Ann Arbor, MI 48109-2121  
shixl@umich.edu

Matthew Bonner  
Department of EECS  
University of Michigan  
Ann Arbor, MI 48109-2121  
mabonner@umich.edu

Lada Adamic  
School of Information  
University of Michigan  
Ann Arbor, MI 48109-1107  
ladamic@umich.edu

Anna C. Gilbert  
Department of Mathematics  
University of Michigan  
Ann Arbor, MI 48109-1043  
annacg@umich.edu

## ABSTRACT

Online networks occupy an increasingly larger position in how we acquire information, how we communicate with one another, and how we disseminate information. Frequently, small sets of vertices dominate various graph and statistical properties of these networks and, because of this, they are relevant for structural analysis and efficient algorithms and engineering. For the web overall, and specifically for social linking in blogs and instant messaging, we provide a principled, rigorous study of the properties, the construction, and the utilization of subsets of special vertices in large online networks. We show that graph synopses defined by the importance of vertices provide small, relatively accurate portraits, independent of the importance measure, of the larger underlying graphs and of the important vertices. Furthermore, they can be computed relatively efficiently.

## Categories and Subject Descriptors

E.1 [Data Structures]: Graphs and networks; G 2.2 [Graph Theory]: Graph algorithms; E.4 [Coding and information theory]: Data compaction and compression

## General Terms

Measurement, Algorithms, Theory

## Keywords

Graph synopsis, graph compression

## 1. INTRODUCTION

Networks play a crucial role in how we acquire information, how we convey information to one another, and how we interact with other people. On the World Wide Web, we do

this through generating and linking content. To study the flow of information, to optimize engineering systems, to design efficient algorithms [5, 13, 14], and to investigate social structure and interaction, we study the statistical and graph properties of *entire* networks, including such features as degree distributions, connectivity, diameter, clustering properties, and evolution of such networks [4, 6]. For a variety of online networks, small subsets of vertices are relevant for efficient algorithms and dominate various graph and statistical properties. Frequently, these smaller subsets or *graph synopses* are easier to study and to understand. One might be interested in whether relationships among web pages can be described without crawling the whole web graph and might be inferred from a small set of vertices. We might also study the “communication” among the most influential political blogs [2] and determine whether information flows directly among them or through intermediate blogs. Despite these examples, there is little principled study of the properties, the construction, and the utilization of subsets of special vertices or edges in large real networks. Such a study is challenging because it is hard to define precisely what is meant by a small version of the graph. Also, it is difficult to evaluate the quality of a compressed graph.

We would like a simple, principled approach to graph synopsis for a number of reasons. First, there are a number of online networks in which a synopsis of the graph is sufficient to capture the relevant information we seek. For example, rather than continuously tracking millions of blogs, one may use occasional snapshots of the blogosphere to construct a subgraph of the most “important blogs” according to a desired measure, and crawl, query, and analyze this smaller synopsis. The synopsis will allow us to capture predominant features of the much larger underlying graph, but, due to its small size, can be stored much more efficiently and even distributed and replicated amongst a number of resource-constrained computers which themselves can execute queries on the content and links.

To build a principled approach to graph synopsis, we start with the definition of predominant vertices and define a precise construction of a graph synopsis from these. Typically, the subset of vertices which capture the graph features are those which are “important.” Furthermore, the importance of these vertices is highly skewed—only few of them are of great importance and the majority are less important. These

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.

HT'08, June 19-21, 2008, Pittsburgh, Pennsylvania, USA.  
Copyright 2008 ACM 978-1-59593-985-2/08/06 ...\$5.00.

vertices and subgraphs have been studied extensively in on-line networks [32, 7], but not with the idea of using them for graph synopses. Following much of this work, we choose four standard definitions of importance: degree, betweenness, closeness and PageRank. We demonstrate empirically for a number of representative online networks that these subsets of vertices do not depend highly on the choice of importance measure. Next, we show that it is possible to glean accurate information about the communication, relationship, and flow of information on the original graph and among the top vertices simply from a subgraph constructed from the important vertices. Furthermore, these properties are consistent, regardless of the importance measure we use, and are appropriate for efficient algorithm design and information management.

There are previous studies about compressing web graphs for space-efficient data storage and transfer [30, 3], using a subgraph to represent the original large graph (the graph sampling problem) [17, 15], mining a subgraph for visualization of the original graph [10, 31], placing sensors to detect information flow [18], constructing a synopsis by projecting queries [16], and quantifying the extent to which important vertices hold online social networks together [21]. However, this work has focused on keeping or representing the properties of the original networks; i.e. studying the entire networks. We study the more fundamental properties of the subgraphs of important vertices themselves.

We give a clear, precise definition of the algorithmic problem of *vertex-importance graph synopsis* in section 2 and discuss the computational hardness of this problem in section 4. We show in sections 3 and 4 that most online networks are far from the worst-case graph; they exhibit features (e.g., power-law degree distribution, short average diameter, and high clustering) that allow us to efficiently compute a graph synopsis. Moreover, we tie properties of the subgraphs to measures, such as assortativity, of the original networks. Finally, in section 5, we match the empirical observations to analytical results.

## 2. PRELIMINARIES

### 2.1 Importance measures

The definitions of *importance* or *prominence* on vertices vary significantly depending on the specific network and application. Most such measures describe the topological location of the vertices. We choose four of the most commonly used measures in various applications as our objects of study: *degree*, *betweenness*, *closeness*, *PageRank*.

Let the graph  $G(V, E)$  have  $|V| = n$  vertices, the four importance values defined on vertices  $v_i$  are listed below:

1. **Degree**  $D(v_i)$ : a measure of how many vertices in  $G$  are connected to  $v_i$  directly. If  $G$  is a undirected graph, then  $D(v_i)$  is the number of undirected edges incident to  $v_i$ ; if  $G$  is a directed graph, then  $D(v_i)$  is the sum of indegree and outdegree of  $v_i$ , where indegree is a count of the number of directed edges to the vertex, and outdegree is the number of directed edges from that vertex to others. Degree reflects a local property of the vertices in the graph.
2. **Betweenness**  $B(v_i)$ : a measure of how many pairs of vertices go through  $v_i$  in order to connect through

shortest paths in  $G$ :

$$B(v_i) = \sum_{j < k} g_{jk}(v_i) / g_{jk}$$

where  $g_{jk}$  is the number of shortest paths linking vertices  $j$  and  $k$ ; and  $g_{jk}(v_i)$  is subset of those paths that contain vertex  $v_i$ . For a directed graph  $G$ , the shortest paths are directed shortest paths. Betweenness reflects a global property of the vertices in the graph.

3. **Closeness**  $C(v_i)$ : a measure of the distances from all other vertices in  $G$  to vertex  $v_i$ :

$$C(v_i) = \left[ \sum_{j \neq i} d(v_i, v_j) \right]^{-1}$$

where  $d(v_i, v_j)$  is the distance between  $v_j$  and  $v_i$ . Intuitively, closeness means that vertices that are in the “middle” of the network are important. For a directed graph  $G$ , the closeness of a vertex could be computed in three ways: all directed paths *to* the vertex, all directed paths *from* the vertex, and all paths regardless of direction. In our work we use this third version, effectively treating the graph as undirected.

4. **PageRank**: a variant of the Eigenvector centrality measure and assigns greater importance to vertices that are themselves neighbors of important vertices [25].

### 2.2 Description of network datasets

We chose our network data sets to be representative of web and online social network data for which one might be interested in examining the properties of important vertices and their graph synopsis. We complement three empirical data sets with analysis of Erdős-Renyi (ER) random graphs, in order to discern interesting features in real world graphs from patterns that may arise by chance. For directed and undirected graphs, we measure the properties of the directed or undirected versions respectively, restricting ourselves to the largest weakly connected component.

	Erdős-Renyi	BuddyZoo	TREC	Web
Vertices	10,000	135,131	29,690	152,171
Edges	49,935	803,200	195,940	1,686,541
ASP	4.26	5.96	3.72	3.48
Directed	False	False	True	True

**Table 1: The average shortest path (ASP) and other characteristics of the largest components of the graphs.**

**Erdős-Renyi random graph.** An Erdős-Renyi random graph is a prototypical random graph with each pair of vertices having an equal probability  $p$  of being joined by an edge. In our model, we set the number of vertices  $|V| = 10000$  and choose  $p = \frac{1}{1000}$ , so the average degree is  $\langle d \rangle = p \times |V| = 10$ .

**Budyzo dataset.** The first real-world network we consider is derived from the website [buddyzo.com](http://buddyzo.com). The system, no longer active, allowed users to submit their AOL Instant Messenger (AIM) buddy lists to compare with others. By treating each registered user as a node and their Buddy List as a series of edges to other nodes, a graph is

formed. Our anonymized snapshot of the data from 2004 includes 140,181 registered users [12]. In this paper, we keep only reciprocal ties (74.7% of the total edges), producing an undirected graph.

**TREC.** The second real-world graph considered is a network of blog connections, the TREC (Text REtrieval Conference) Blog-Track 2006 dataset [20]. It is a crawl of 100,649 RSS and Atom feeds collected over 11 weeks, from December 6, 2005 to February 21, 2006. In our experiments, we removed duplicate feeds and feeds without a homepage or permalinks. We also removed over 300 Technorati tags, which appear to be blogs, but are in fact automatically generated from tagged posts, and so are not true indicators of social linking. The TREC dataset contains hyperlinks of various forms, including blogrolls, comments, trackbacks, etc. There are 198,141 blog-to-blog hyperlinks in total, and 33,385 blogs having at least one such link.

**Web graph dataset.** The web graph data set was collected in 1998 by Alexa<sup>1</sup> and has previously been analyzed as part of the “Web in a box” project at the Xerox Palo Alto Research Center [1]. Since the snapshot was collected such a long time ago, it contains *only* 50 million pages and 259,794 websites. This “small” size allows us to comprehensively analyze the web graph. We construct a directed graph where Site A has a directed edge to site B if any of the pages within A point to any page within site B.

Due to the similarity of results for the recent blog datasets and the decade old website-level data set, we expect our results to be applicable to larger, more current webcrawls.

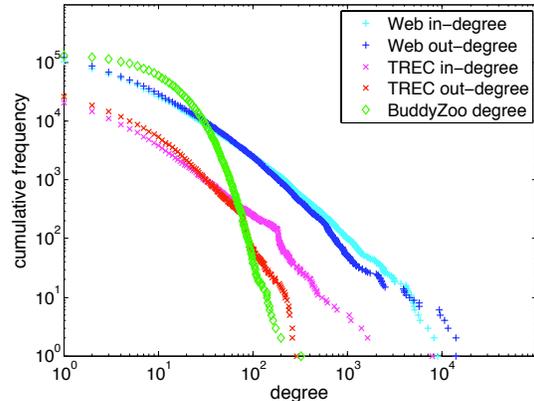
### 3. IMPORTANT VERTICES

In this section, we examine the graph synopsis consisting of important vertices in the network. First, we describe some properties of the entire networks. Second, we analyze the subgraphs induced by important vertices. Finally, we compare some properties of the important vertices in the subgraphs and the entire networks.

#### 3.1 Network properties and important vertices

**Degree distributions.** We plot the cumulative degree distributions of three real online networks in Figure 1. We treat the Web and TREC networks as directed graphs and plot the distributions of their in-degrees and out-degrees and we treat BuddyZoo as an undirected graph. By fitting the distributions of in-degree of Web and TREC with power-law distributions, we get their power-law exponents, which are 2.47 and 2.16 respectively. Moreover, we can see that the degree distribution of BuddyZoo has a very sharp drop off at the tail, which is observed in many social networks, e.g. co-authorship networks [22]. This places blog links, a form of social linking, somewhere between navigational/informational general linking on the Web and the reciprocal, communicative linking of a social network. The distributions of out-degree of Web and TREC show mild deviations from power laws, consistent with other web measurements [28] and might be due to the limitation of the data sampling [29].

**Correlation of importance values of different measures.** Before examining the important vertices in the networks, we look at the relationships of importance measures in different networks. Table 2 shows that all of the im-



**Figure 1: The degree distributions of online networks of BuddyZoo data, TREC blog data and Web data.**

portance measures are positively correlated in all four networks. The two undirected graphs, Erdős-Renyi and BuddyZoo, have more highly correlated importance measures. Perhaps the directed edges of the other graphs add complexity to centrality measures. Furthermore, we see that for all of the networks, *degree*, *betweenness* and *PageRank* have higher correlation than *closeness*. Thus, we see that there are various ways of defining importance in the networks and the most *central* vertices according to different centrality measures share overlap significantly.

Correlation	Erdős-Renyi	BuddyZoo	TREC	Web
Deg, Bet	0.9920	0.8137	0.7872	0.6178
Deg, Clo	0.9474	0.7849	0.3835	0.7869
Deg, PR	0.9952	0.9486	0.7058	0.5175
Bet, Clo	0.9673	0.7541	0.3120	0.4709
Bet, PR	0.9823	0.8439	0.7439	0.6757
Clo, PR	0.9154	0.6418	0.1086	0.3253

**Table 2: Spearman correlations between importance measures of vertices. All the  $p$ -values of the correlations are  $< 0.0001$ .**

**Assortativity.** The concept of *assortativity* or assortative mixing is defined as the preference of the vertices in a network to have edges with others that are similar. Here, we will focus on similarity with regard to centrality. We choose to measure the average value  $\langle k \rangle$  of the neighbors of vertices of importance value  $k$ , i.e.  $\langle k \rangle_{neigh}(k) = \sum_{k'} k' P(k'|k)$ , where  $k$  is determined by each of the four different importance measures [27]. From the change of  $\langle k \rangle_{neigh}(k)$  as  $k$  increases, we deduce the network’s assortativity for this particular valuation. When the overall slope of  $\langle k \rangle_{neigh}(k)$  is positive, the network is assortative; if the overall slope is negative, then it is disassortative. Otherwise, the network is neutral (e.g. the assortativity of degree of Erdős-Renyi random graphs).

In Figure 2, we can see that all four networks are consistently assortative with the importance measure of closeness. This confirms our intuition—the neighbors of the vertices with high closeness also have high closeness. The other three importance measures consistently show that the

<sup>1</sup>www.alexa.com

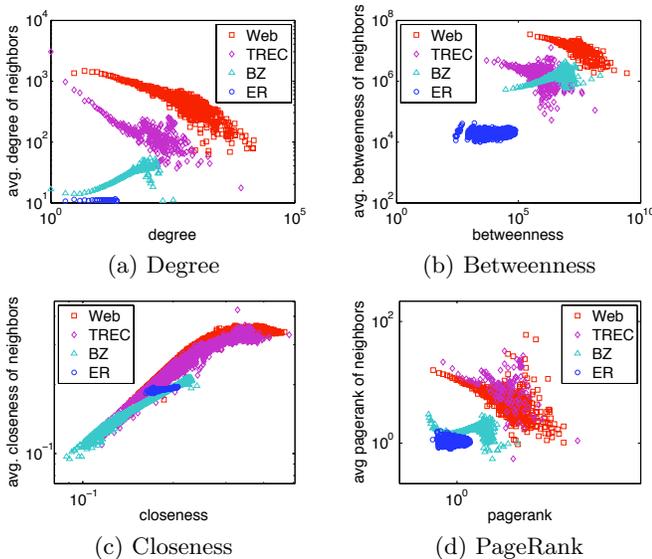


Figure 2: The slopes of the distributions of  $\langle k \rangle_{\text{neigh}}$  show the assortativities.

Erdős-Renyi random graph is a neutral graph, that BuddyZoo, similar to other social networks[23], is assortative, and that the Web and TREC blog networks are mildly disassortative. We'll see in Section 5 that this result does not mean that important blogs avoid linking to other important blogs. Rather, there is such a large skew in the linking behavior to important blogs, that one would expect at random for them to already be linking to one another very frequently.

### 3.2 Important vertices in their subgraphs

In this section, we discuss important vertices and the subgraphs induced by these vertices. Such analysis helps us to discover the information hidden behind the important vertices in the real online networks, and how we can utilize them for graph synopsis. We do not fix a specific threshold for inclusion of important vertices in the subgraph, as this may vary by application. Rather, in our study what occurs as we allow the absolute number of important vertices  $m$  vary, as long as  $m \ll n$ , where  $n$  is the number of vertices in the original network.

Figure 3 shows the subgraphs induced by the four importance measures in BuddyZoo and the highest degree vertices in the other three networks. These subgraphs may be markedly different for different measures of importance, even within the same graph, in spite of high correlation in importance measures among vertices. They may also vary significantly between graphs, even for the same importance measure. There are several explanations of this behavior. Given the high assortativity of the closeness measure, we are unsurprised to find that individuals of high closeness are closely connected in the BuddyZoo graph. Buddyzoo also has individuals of high degree, but there were limits imposed both by AOL and individuals' own bandwidth, and so the large connected component among high degree vertices does not contain all such vertices. On the other hand, the highest degree vertices in both the blog and web datasets have such high degree that they tend to form a single connected component.

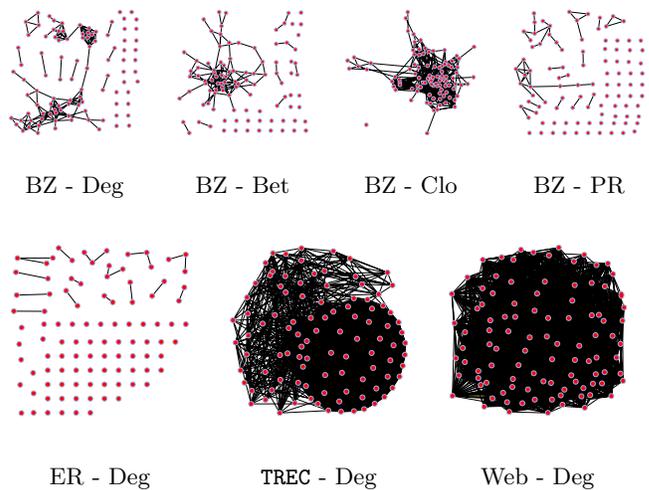


Figure 3: In the top row are subgraphs induced by the top 100 important vertices of BuddyZoo for all four importance measures, while in the bottom row are subgraphs induced by the 100 highest degree vertices in the other three networks.

**Connectivity.** The first question we address is whether the connectivity of important vertices depends on other, less important, vertices or whether they are already well connected through one another. In the Erdős-Renyi random graph, the size of the largest connected component is given by the solution  $x$  to the equation

$$x = 1 - e^{-\langle k \rangle x}$$

where  $\langle k \rangle$  is the average degree of the graph. The solution to this equation, shown as a dotted red line in Figure 4(a), represents the change in size of the largest connected component of the subgraphs induced by picking vertices randomly from the Erdős-Renyi random graph. When we choose vertices according to importance instead, the subgraphs have significantly better connectivities, with the largest connected component occupying 96.5% of the subgraph once the subgraph contains over 15% of all vertices in the graph (i.e., 1,500 important vertices vs. 10,000 total vertices).

Moreover, from Figure 4, we see that the important vertices are even more highly connected in the real networks. No matter which network and which importance measure, all of the curves of the connectivity of important vertices are almost monotonically increasing. For BuddyZoo, more than 95% of the important vertices of highest degree, betweenness or closeness are in the largest connected component when they comprise just 1% percent of all vertices in the network (i.e., 1,500 important vertices vs. 135,131 total vertices). In addition, more than 95% of the 10,000 highest PageRank vertices are in the largest connected component. For both of the two directed networks, the TREC blog network and the network of websites, the most important vertices are very well connected ( $> 99.5\%$ ) even when their numbers are very small ( $< 0.05\%$  of all the vertices in the networks). Note that this very high level of connectivity is in spite of the dissortative nature of the TREC and website networks with respect to degree, betweenness and PageRank, where

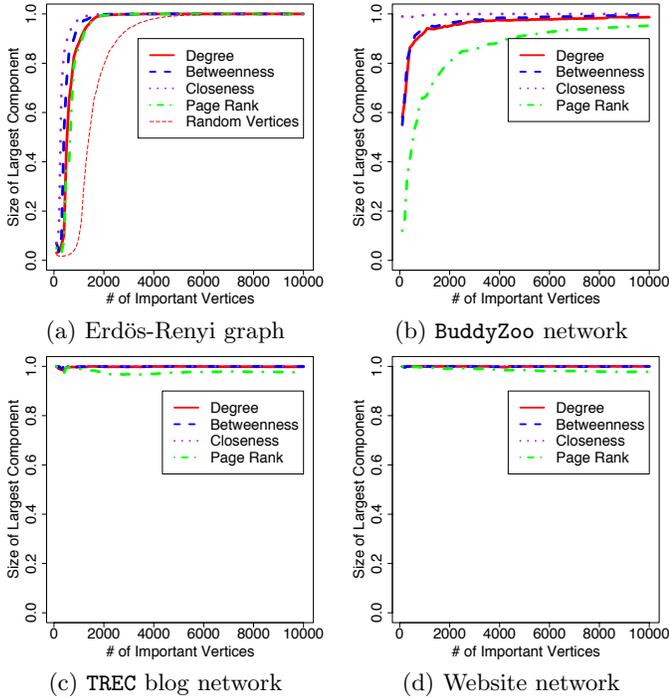


Figure 4: The sizes of largest connected component of the sub-networks of important vertices in Erdős-Renyi random graph and three real online networks.

important vertices tend to connect to less central vertices. We can reconcile the two by observing that the important vertices are already interconnected, so the negative assortativity comes from highly connected vertices being connected to lower degree vertices simply because they already have so many connections and there is only a small percentage of vertices of similarly high degree [26].

**Density.** The previous observations tell us that the connectedness of important vertices is high even when we omit all other vertices in the original graph and even when they comprise a very small fraction of the entire network. Next, we examine just how dense their connections are. In Figure 5, we show the relationships between the number of edges incident on important vertices and the number of important vertices.

Figure 5 (a) shows that for an Erdős-Renyi graph, the important vertices according to all four measures have a higher average degree in the subgraph than randomly chosen vertices (red dashed line), but this average degree is lower than the average degree in the complete graph (black dashed line). The density of the graph reaches a maximum when all of the vertices in the graph are included. Moreover, from the direction of the curves, we can see that the number of edges  $e$  increases super-linearly with the number of important vertices  $n$ , i.e.  $\Theta(n) < e < \Theta(n^2)$ .

However, Figures 5(c) and (d) reveal the opposite behavior for networks with highly skewed degree distributions (TREC and Web). The curves of each network do not overlap as much, and the average degree of the important vertices in the subgraph is higher than the average degree in the original network. This indicates that rather than being sparser, as

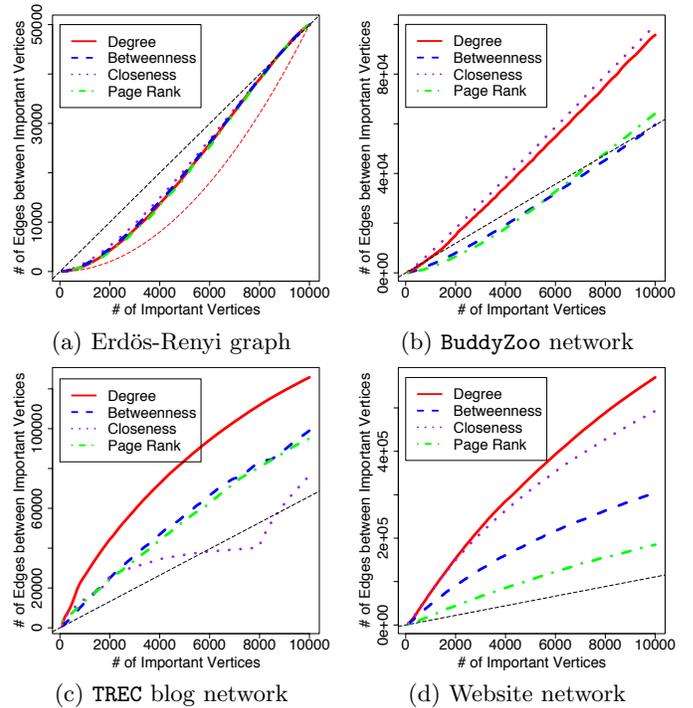


Figure 5: The growth of numbers of edges between important vertices. The slope of the black dash line in each plot is the ratio of the number of edges v.s. the number of vertices in the entire network.

was the case for the Erdős-Renyi subgraphs, the subgraphs of important vertices in real world online networks are actually *denser*. Finally, for the BuddyZoo network (Figure 5 (b)), which is assortative, but not power-law in degree, we see a mix of trends. Subgraphs of vertices with high betweenness and PageRank tend to be a bit sparser than the complete network, but the most important vertices according to degree and closeness are more densely connected (this is also apparent in the visualizations in Figure 3).

In examining these real online networks, we see that although the densities of connection among important vertices vary considerably in different networks with different importance measures, in general, they are significantly denser than random vertices in the Erdős-Renyi random graph.

### 3.3 Original vs. subgraph properties

**Distance.** In Section 3.2 we saw that even without any additional vertices from the original graph, the subgraphs of important vertices in the three online networks are already well connected. Next we examine the second property that we want to preserve for our graph synopsis problem: the average shortest paths (ASP) between reachable pairs of important vertices.

Figure 6 shows the comparison curves of ASPs of important vertices in their induced subgraphs and in the original networks. In the Erdős-Renyi random graph, the ASP between important vertices is on average shorter than the ASP for the entire network (the dotted baseline). But in their induced subgraphs there are significantly more hops between them on average, which indicates that important

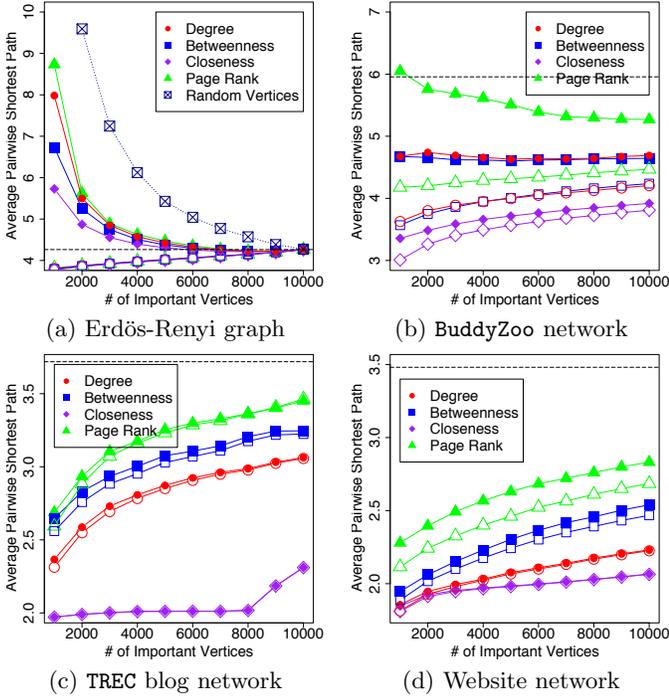


Figure 6: The ASP of: all vertices in the entire networks (black dashed line), important vertices in the the subgraphs (solid points), important vertices in the entire networks (hollow points).

vertices in random networks are not closely connected, and their shortest paths route through non-important vertices. Nevertheless, subgraphs of important vertices in ER graphs are somewhat better connected than subgraphs of randomly selected vertices.

In contrast to the Erdős-Renyi random graph, all three real online networks consistently show that the ASPs between important vertices are much shorter than the average shortest paths of the entire network; and almost all of them are increasing as more important vertices are added in. What is more, by comparing the ASPs of important vertices in the original graphs and in the subgraphs, we see that their values are extremely close in most cases, especially for the TREC and Web data, e.g. the solid and hollow purple points (ASP of vertices of highest closeness) are almost exactly overlapped. This indicates that important blogs are most efficiently connected through other important blogs.

**Relative importance.** In addition to the connectedness of important vertices, we are also interested in their relative ranking: if we only keep the important vertices and the edges among them, how would the vertices rank in the new subgraph with the same importance measure? In order to answer this question, we generate subgraphs of different sizes for all networks. We then compute the importance of the vertices in the subgraphs according to the same importance measure used to select them. Finally, we compute the Pearson correlation of the importance values of those vertices in the original graph and in the subgraph.

Figure 7 shows that the correlations are all much higher for the real-world online networks than the Erdős-Renyi ran-

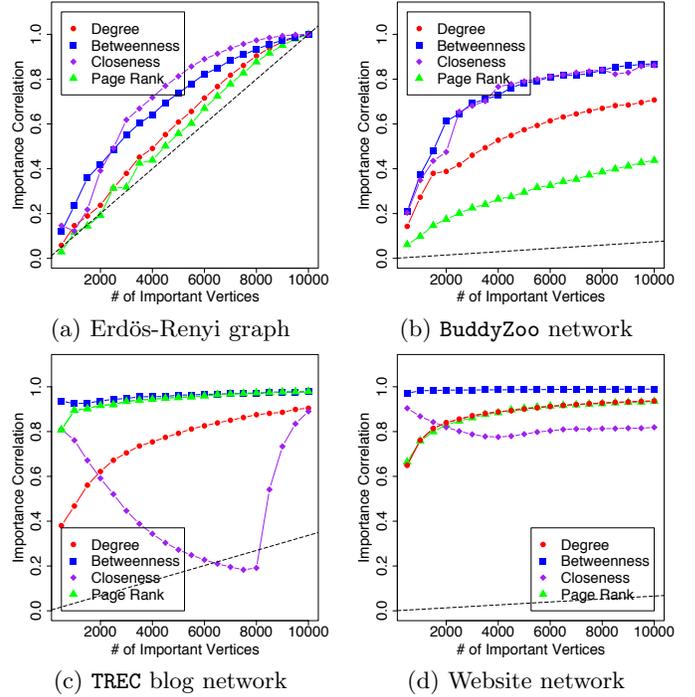


Figure 7: Pearson correlations of importance values of vertices in subgraphs and original graphs. The black dashed lines are the base lines starting from 0 when the number of vertices is 0; and ending at 1 when all the vertices in the networks are included.

dom graph, and that this is especially true for the Web and TREC data. The abnormality of closeness in TREC may be due to the blog aggregators and splogs. There are one to two thousand blogs whose only incoming link is from centrally positioned (in the network) blog aggregators. This boosts the closeness score of the unimportant blogs, creating the abnormality in Figure 7.

The high correlations of the online networks tell us that the ranking of importance in the subgraphs of important vertices is highly consistent with their ranking in the original graphs. This suggests that, e.g. it may not be necessary to crawl all blogs to get an accurate ranking of the most important blogs. Rather, the links among the top blogs themselves may already provide fairly close approximate rankings.

### 3.4 Summary

After studying the important vertices and their induced subgraphs, we can make two overall observations about the four networks: (i) Different importance measures yield subgraphs of varying density and topology as is evident in Figure 3. (ii) However, in spite of these differences, “important vertices” in the online networks have some properties that agree with each other, which are essential for the graph synopsis we are looking at: they connect to each other more directly than average; their distances to each other are closer than random vertices; and their relative ranks are positively correlated to their importance ranks in the original networks. Thus, we know that in the real online networks, in contrast to random graph model, the subgraphs induced by

the important vertices tend to preserve information about the relationships among important vertices, and we can use the subgraphs to study the properties of important vertices in the original graphs.

## 4. COMPRESSION WITH GUARANTEES

While retaining only the important vertices may be sufficient to capture most of the relationships among them in real-world networks, in general we have no guarantee that these induced subgraphs preserve any properties at all (whether of the important vertices or of the original graph). We cannot even guarantee the most basic property of connectivity of the important vertices. In this section, we rigorously define the graph compression problem, analyze the computational complexity of two heuristic algorithms, and discuss the trade-offs of these approaches.

### 4.1 Hardness of compression with guarantees

We define the BASIC GRAPH COMPRESSION PROBLEM as follows: In a connected unweighted graph  $G(V, E)$ , every vertex is assigned an importance value. Taking the original graph  $G(V, E)$  and the set of vertices  $S$  with largest importance values as inputs, find the minimal set of additional vertices  $\nu$ , which form a connected subgraph  $G'(V', E')$ , where  $V' = S + \nu$  and  $V' \subseteq V, E' \subseteq E$ .

We recall the NETWORK STEINER TREE PROBLEM which is NP-complete [9]. A heuristic method, the *Minimal Spanning Tree* algorithm gives solutions to this problem with approximation ratio 2 [11]. One can show that BASIC GRAPH COMPRESSION and the NETWORK STEINER TREE PROBLEM are polynomial-time reducible to one another. Thus, BASIC GRAPH COMPRESSION is an NP-complete problem.

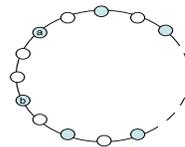
### 4.2 Heuristic algorithms

There are, however, several heuristic algorithms that guarantee the preservation of some properties of the important vertices in the original graph. We detail the KEEPONE and the KEEPALL algorithms [10] next, and note the similar web projection method [16].

**KEEPONE.** Let  $K_1$  be the set of important vertices, the goal is to find the minimal set  $K_2$  such that there is a tree induced by  $K_1 \cup K_2$ . The approximation algorithm is first to build a minimum spanning tree on the complete graph on  $K_1$  where an edge  $(u, v)$  has weight equal to the length of a shortest path from  $u$  to  $v$ . The set  $K_2$  consists of any additional vertices along any “path” edge in the minimum spanning tree. The result is the graph induced by the vertices  $K_1 \cup K_2$ .

The KEEPONE algorithm guarantees the connectivity of the compressed graph, has the same set of additional vertices as the projection method in [16], and only introduces more edges, which means it may have better diameter preservation than the projection method.

Unfortunately, retaining only connectivity may provide a distorted view of the original graph. We see in Figure 8 an example of a graph on  $n$  vertices in which the distance of the original vertices  $a$  and  $b$  is 3 but in the compressed graph built by KEEPONE, their distance is  $n - 3$ . The ratio of the distances is  $\frac{n-3}{3}$  which we can make arbitrarily large by increasing the number of vertices  $n$ . That is, KEEPONE retains connectivity but may drastically distort the distance between some pairs of important vertices. To ameliorate this problem, one can use the KEEPALL algorithm [10] which keeps vertices that lie along a shortest path between any two vertices in  $K_1$ .



**Figure 8:** The distance of important vertices  $a$  and  $b$  in the original graph is 3 and  $n - 3$  in the compressed graph obtained by KEEPONE. The ratio of distances can be made arbitrarily large as  $\lim_{n \rightarrow \infty} \frac{n-3}{3} = \infty$ .

### 4.3 Empirical evaluation and trade-offs

While Figure 8 shows that the worst case distance preservation of KEEPONE may be arbitrarily bad, real-world networks are far from the worst case. Furthermore, the KEEPONE and KEEPALL algorithms illustrate that there are some tradeoffs we may make in compressing real-world graphs—we can maintain distances at the cost of keeping a few additional vertices. To explore these empirical tradeoffs, we apply both the KEEPONE and KEEPALL algorithms to three networks. Table 3 shows these results. Since the results with the Web data are very similar to TREC, we do not list them here for conciseness. From the table, we can see that if we insist on preserving the pairwise shortest paths of all important vertices, we must include many more additional vertices (thus increasing the size of our synopsis). Furthermore, we must do so even though the average pairwise shortest paths in the subgraph of just the important vertices is already close to that of the original graph. Note that we increase the size of the synopsis by fewer than 100 additional vertices when we preserve connectivity (with KEEPONE), but we need over 3000 additional vertices when we also insist on preserving distances. In short, while the problem of preserving connectivity in graph compression is NP-complete, heuristic algorithms such as KEEPONE can preserve connectivity with a lower cost, while preserving the distances demands quite more. In this sense, we can also see that the short pairwise shortest paths of important vertices in their subgraphs and their original graphs is a special and important property of the online networks we study.

## 5. ANALYTICAL DISCUSSIONS

In this section we present the expected density of subgraphs of random graphs with varying degree distributions, in order to contrast these expected values with the empirically observed measurements. We limit ourselves to vertex degree as the sole importance measure and assume that the graphs are random aside from the degree distribution, which we specify. We then obtain the density of the subgraph by deriving the probability that an edge in the original graph lies between two vertices in the subgraph.

First, we find the degree  $k_i$  of the least important vertex among the set of top  $i$  most important vertices. We do so by calculating the expected number of vertices of degree at least  $k_i$  in a network with  $n = |V|$  vertices. Furthermore, we assume that the expected number is actually equal to  $i$  so that

$$i = n \cdot P(k_i).$$

Because we are given the cdf  $P(k)$  explicitly for Erdős-Renyi and power law random graphs, we can solve the pre-

Subgraph	Add vts	LC	Avg PSP	Subgraph	Add vts	LC	Avg PSP	Subgraph	Add vts	LC	Avg PSP
Erdős-Renyi				BuddyZoo				TREC			
Sub-Deg100	0	NA	NA	Sub-Deg100	0	0.58	NA	Sub-Deg100	0	1	1.636
KO-Deg100	80	1	14.526	KO-Deg100	33	1	9.440	KO-Deg100	0	1	1.636
KA-Deg100	3222	1	3.649	KA-Deg100	2199	1	3.233	KA-Deg100	34	1	1.609
Sub-Bet100	0	NA	NA	Sub-Bet100	0	0.55	NA	Sub-Bet100	0	1	2.085
KO-Bet100	68	1	15.497	KO-Bet100	35	1	16.087	KO-Bet100	0	1	2.085
KA-Bet100	3185	1	3.633	KA-Bet100	2376	1	3.171	KA-Bet100	216	1	1.994
Sub-Clo100	0	NA	NA	Sub-Clo100	0	0.99	2.599	Sub-Clo100	0	1	1.716
KO-Clo100	62	1	11.474	KO-Clo100	1	1	2.624	KO-Clo100	0	1	1.716
KA-Clo100	3000	1	3.604	KA-Clo100	531	1	2.324	KA-Clo100	0	1	1.716
Sub-PR100	0	NA	NA	Sub-PR100	0	0.12	NA	Sub-PR100	0	1	1.298
KO-PR100	87	1	15.404	KO-PR100	75	1	11.517	KO-PR100	0	1	1.298
KA-PR100	3338	1	3.672	KA-PR100	3978	1	3.880	KA-PR100	36	1	1.294

**Table 3: Comparison of the properties of subgraphs generated by different methods with important vertices in Erdős-Renyi random graph, BuddyZoo and TREC. Sub-ImportanceMeasure100 is the subgraph induced by top 100 important vertices only; KO- is the subgraph generated by `KEEPONE`; KA- is the subgraph generated by `KEEPALL`. LC is the fraction of important vertices in the large component of the subgraph. Avg PSP is the average pairwise shortest path length in the subgraph.**

vious equation for  $k_i$  and, after doing so, we find the probability that an edge is incident to a single important vertex,  $e \rightarrow V_i$ , given by

$$\mathbb{P}(e \rightarrow V_i) = \frac{1}{|E|} \int_{k_i}^n k \cdot p(k) dk$$

where  $p(k)$  is the pdf of the degree distribution. Using independence of the edges, we find that the number of edges within the subgraph of important vertices is simply

$$|E_i| = |E| \cdot \mathbb{P}(e \rightarrow V_i)^2.$$

## 5.1 Erdős-Renyi graphs

In an Erdős-Renyi random graph, the degrees are distributed according to a Poisson distribution where the probability of a vertex having degree larger than the mean decreases exponentially. As a result, even when selecting the highest degree nodes, their degree will be within an order of magnitude of the average degree  $z = \langle k \rangle$  of the network.

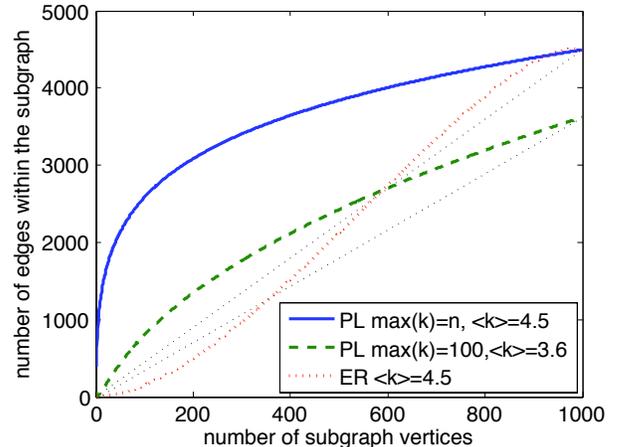
In Figure 9, we show the number of edges in the subgraph of an Erdős-Renyi graph, using the normal distribution with mean  $z$  and standard deviation  $\sigma = \sqrt{z/n} * (1 - z/n)$ , is

$$i = \frac{1}{2} \left( 1 + \operatorname{erf} \left( \frac{k_i - z}{\sigma\sqrt{2}} \right) \right).$$

We see that when the number of important vertices is small, the degree within the subgraph is lower than the degree of the original graph. Using well known properties of Erdős-Renyi graphs, we expect that when the average subgraph degree is 1, a giant component will emerge in the subgraph, and further, when the average degree is  $\log(n)$ , the subgraph will be path connected. This is consistent with the set of connectivity and density measurements on simulated Erdős-Renyi graphs in Section 3.2.

## 5.2 Power law graphs

We expect different behavior in power law graphs, where high degree vertices are so well connected, that they will



**Figure 9: The number of edges between important vertices, where importance is measured by degree, in three networks: 1) power law network with  $\alpha = 2.2$ ,  $n = 1000$ , 2) Erdős-Renyi graph with the same average degree, and 3) power-law graph with the same exponent but a cutoff at  $k = 100$ . Two dotted lines show what the number of edges would be if the average degree in the subgraph were equal to the average degree in the original network.**

naturally connect not only to a large portion of the network, but to one another as well. For example, in a power-law graph with exponent  $\alpha$  and no cutoff on the degree<sup>2</sup>, one vertex on average is expected to have degree  $N^{1/(\alpha-1)}$  [24].

<sup>2</sup>a cutoff may be imposed such that  $P(k) \sim k^{-\alpha}$  for  $k < \max(k)$  and 0 otherwise

For  $\alpha = 2$ , this means that one expects one node to be connected to majority of the other nodes.

In selecting high degree nodes in a power law graph, we are selecting nodes that are likely to be connected to each other by virtue of the fact that so many edges are incident on them. The number of vertices with degree  $k_i$  or greater is given by

$$i = n P(k \geq k_i) = \frac{n}{k_i^{\alpha-1}}$$

Solving for  $k_i$ , we have that the degree of the  $i^{th}$  most important vertex is  $k_i = (\frac{n}{i})^{\frac{1}{\alpha-1}}$ . Next, we want to find out what proportion of the edges are incident on the  $i$  most important vertices. For this we have

$$P_e(i) = P(e \in e_i) = \frac{\int_{k_i}^n kp(k)dk}{\int_1^n kp(k)dk} \quad (1)$$

$$= \frac{k_i^{2-\alpha} - n^{2-\alpha}}{1 - n^{2-\alpha}} = \frac{(\frac{n}{i})^{\frac{2-\alpha}{\alpha-1}} - n^{2-\alpha}}{1 - n^{2-\alpha}} \quad (2)$$

Figure 9 shows that the average degree in the subgraphs of important vertices is actually *higher* than in the original graph. We repeat the analysis using a degree distribution cutoff  $\max(k)$  that is lower than the total number of nodes  $n$ . This cutoff not only disallows very high degree vertices, but also lowers the average degree in the original subgraph. When the cutoff is introduced, the subgraph still maintains a higher average degree than the original graph, but the difference is less pronounced.

Note the similarity with Figure 5, showing the number of edges in the subgraph for the TREC and Web data sets, both of which are power law in nature (although directed). In both the analytical and empirical subgraphs, the average degree is higher than it is for the entire graph. We should mention that for exponents  $\alpha \sim 2$  and very small  $i$ , Equation 2 would yield a higher average degree than there are important vertices to connect to. This is in fact a known property of random power law graphs, where simply fixing the degree of a vertex and allowing it to satisfy this degree by forming edges at random would create a non-vanishing frequency of multiple edges between highly connected vertices. If one disallows multiple edges, the networks become mildly disassortative, consistent with our empirical observations.

## 6. RELATED WORK

In this section, we examine the graph sampling problem and the rich-club phenomenon. Both of them have some similarities with our problem: the former also studies how to get “good” subgraphs given large massive networks; and the later focuses on the set of “important vertices”. However, they are still different from our problem in various aspects. In graph sampling, one aims to devise a sampling method, e.g. random vertex or edge selection, snowball sampling, the sketching-based sampling [19] etc., in order to be able to infer the properties of the original graph from the much smaller sampled graph [15, 17]. In contrast, our work constructs subgraphs of predetermined important vertices, not for the purpose of deducing properties of the original graph, but in order to infer the underlying relationships amongst the important vertices themselves.

In the “rich-club phenomenon”, vertices with high degree tend to connect together tightly, which is true for many social and other types of real networks [32, 7]. While previous

work on the rich-club phenomenon has aimed to determine whether the number of edges between high degree vertices based purely on degree is higher than what one would expect at random, our study extends to other centrality measures, and describes essential properties of the subgraphs themselves, such as connectivity, shortest paths, and preserving rank orderings of importance. A related analysis of highly interconnected sub-structures in networks is that of  $k$ -cores, subgraphs of vertices where each vertex has at least  $k$  connections within the subgraph [8]. An interesting direction for future work would be to repeat our analysis of the properties of the subgraph and original graph, using  $k$ -core membership as the importance measure for vertex selection.

## 7. CONCLUSION

In this paper, we propose a new approach to analyzing and studying large online networks, *vertex-importance graph synopsis*. Given a set of important vertices, we extract a much smaller subgraph from the original network, containing those important vertices. We attempt to place this process on a rigorous footing and show that even simple versions of the graph compression problem are hard (but that there are reasonable heuristic algorithms). Unlike previous methods which evaluated the fidelity of the “graph abstract,” this approach utilizes the subsets of important vertices and edges and the information they could provide in large networks. We argue that they can make information access and management more efficient in real applications. These observations suggest future work in using graph synopses for information retrieval and information flow detection.

>From our empirical analysis of three real online networks, we find a number of interesting properties. The important vertices are much more closely and densely connected to each other. They also have significantly shorter pairwise paths, which do not heavily depend on the rest of vertices in the networks, (i.e. their pairwise shortest paths in the subgraphs induced by themselves are close to those in the original graphs). Finally, their relative ranks are almost all highly correlated to their ranks in the original networks. Although our experiments show that the properties of vertices of different importance measures in different networks do vary in some ways, the observations stated above are consistent no matter the type of networks (either social or technological), and regardless of the importance measure we choose. Thus, we may use vertex-importance graph synopses as small but accurate representatives of the important vertices in the larger graph (and, sometimes, of the larger graph itself). Furthermore, the real online networks are relatively easy to compress while preserving important graph properties (they do not exhibit the worst-case behavior of our theoretical analysis).

In addition to empirical studies, we use analytical discussions to show how these properties of important vertices in online networks differ from random graph models. What is more, we also use heuristic algorithms to measure the complexities and trade-offs of requiring some properties of the real networks to be guaranteed in the compressed graphs.

## 8. ACKNOWLEDGMENTS

We thank Martin Strauss for many helpful discussions. We are also grateful to Adam D’Angelo for sharing the BuddyZoo dataset with us.

## 9. REFERENCES

- [1] L. Adamic. The Small World Web. *Proceedings of ECDL*, 99:443–452, 1999.
- [2] L. A. Adamic and N. Glance. The political blogosphere and the 2004 u.s. election: Divided they blog. In *Proceedings of LinkKDD-2005*, 2005.
- [3] M. Adler and M. Mitzenmacher. Towards compressing web graphs. In *DCC '01: Proceedings of the Data Compression Conference (DCC '01)*, page 203, Washington, DC, USA, 2001. IEEE Computer Society.
- [4] Y.-Y. Ahn, S. Han, H. Kwak, S. Moon, and H. Jeong. Analysis of topological characteristics of huge online social networking services. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 835–844, New York, NY, USA, 2007. ACM Press.
- [5] S. Brin and L. Page. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30(1–7):107–117, 1998.
- [6] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. Graph structure in the web. *Comput. Networks*, 33(1–6):309–320, 2000.
- [7] V. Colizza, A. Flammini, M. A. Serrano, and A. Vespignani. Detecting rich-club ordering in complex networks. *NATURE PHYSICS*, 2:110, 2006.
- [8] S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes. k-core organization of complex networks. *Physical Review Letters*, 96:040601, 2006.
- [9] M. Garey and D. Johnson. The rectilinear Steiner problem is NP-complete. *SIAM J. Appl. Math.*, 32(4):826–834, 1977.
- [10] A. C. Gilbert and K. Levchenko. Compressing network graphs. In *LinkKDD*, 2004.
- [11] E. N. Gilbert and H. O. Pollak. Steiner minimal trees. *SIAM J. Appl. Math.*, 16(1):1–29, 1968.
- [12] T. Hogg and L. Adamic. Enhancing reputation mechanisms via online social networks. *Proceedings of the 5th ACM conference on Electronic commerce*, pages 236–237, 2004.
- [13] J. M. Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46(5):604–632, 1999.
- [14] J. M. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan, and A. S. Tomkins. The Web as a graph: Measurements, models and methods. *Lecture Notes in Computer Science*, 1627:1–17, 1999.
- [15] S. H. Lee, P.-J. Kim, and H. Jeong. Statistical properties of sampled networks. *Physical Review E*, 73:016102, 2006.
- [16] J. Leskovec, S. Dumais, and E. Horvitz. Web projections: learning from contextual subgraphs of the web. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 471–480, New York, NY, USA, 2007. ACM Press.
- [17] J. Leskovec and C. Faloutsos. Sampling from large graphs. In *KDD '06: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 631–636, New York, NY, USA, 2006. ACM Press.
- [18] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, and N. Glance. Cost-effective outbreak detection in networks. *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 420–429, 2007.
- [19] P. Li, K. Church, and T. Hastie. A sketch-based sampling algorithm on sparse data. Technical report, Department of Statistics, Stanford University, 2006.
- [20] C. Macdonald and I. Ounis. The TREC Blogs06 collection: Creating and analysing a blog test collection. Tech report (dcs), Dept of Computing Science, University of Glasgow, 2006.
- [21] A. Mislove, M. Marcon, K. Gummadi, P. Druschel, and B. Bhattacharjee. Measurement and analysis of online social networks. *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*, pages 29–42, 2007.
- [22] M. Newman. Coauthorship networks and patterns of scientific collaboration, 2004.
- [23] M. E. J. Newman. Assortative mixing in networks. *Phys. Rev. Lett.*, 89(20):208701, Oct 2002.
- [24] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.
- [25] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford Digital Library Technologies Project, 1998.
- [26] J. Park and M. Newman. Origin of degree correlations in the Internet and other networks. *Physical Review E*, 68(2):26112, 2003.
- [27] R. Pastor-Satorras, A. Vazquez, and A. Vespignani. Dynamical and correlation properties of the internet. *Physical Review Letters*, 87:258701, 2001.
- [28] D. Pennock, G. Flake, S. Lawrence, E. Glover, and C. Giles. Winners don't take all: Characterizing the competition for links on the web. *Proceedings of the National Academy of Sciences*, 99(8):5207, 2002.
- [29] X. Shi, B. Tseng, and L. Adamic. Looking at the Blogosphere Topology through Different Lenses. In *Proceedings of the International Conference on Weblogs and Social Media (ICWSM 2007)*, volume 1001, page 48109, 2007.
- [30] T. Suel and J. Yuan. Compressing the graph structure of the web. In *Data Compression Conference*, pages 213–222, 2001.
- [31] A. Y. Wu, M. Garland, and J. Han. Mining scale-free networks using geodesic clustering. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 719–724, New York, NY, USA, 2004. ACM Press.
- [32] S. Zhou and R. J. Mondragon. The Rich-Club Phenomenon In The Internet Topology. *IEEE Commun. Lett.*, 8:180–182, 2004.