



# Limiting the Neighborhood: De-Small-World Network for Outbreak Prevention

Ruoming Jin<sup>1</sup>, Yelong Sheng<sup>1</sup>, Lin Liu<sup>1</sup>, Xue-Wen Chen<sup>2</sup>,  
and NhatHai Phan<sup>3</sup>(✉)

<sup>1</sup> Department of Computer Science, Kent State University, Kent, USA  
{jin,ysheng,lliul}@cs.kent.edu

<sup>2</sup> Department of EECS, The University of Kansas, Lawrence, USA  
xwchen@ku.edu

<sup>3</sup> College of Computing, New Jersey Institute of Technology, Newark, USA  
phan@njit.edu

**Abstract.** In this work, we study a basic and practically important strategy to help prevent and/or delay an outbreak in the context of network: limiting the contact between individuals. In this paper, we introduce the *average neighborhood size* as a new measure for the degree of being small-world and utilize it to formally define the *de-small-world* network problem. We also prove the NP-hardness of the general reachable pair cut problem and propose a greedy edge betweenness based approach as the benchmark in selecting the candidate edges for solving our problem. Furthermore, we transform the de-small-world network problem as an OR-AND Boolean function maximization problem, which is also an NP-hardness problem. In addition, we develop a numerical relaxation approach to solve the Boolean function maximization and the de-small-world problem. Also, we introduce the *short-betweenness*, which measures the edge importance in terms of all short paths with distance no greater than a certain threshold, and utilize it to speed up our numerical relaxation approach. The experimental evaluation demonstrates the effectiveness and efficiency of our approaches.

## 1 Introduction

The interconnected network structure has been recognized to play a pivotal role in many complex systems, ranging from natural (cellular system), to man-made (Internet), to the social and economical systems. Many of these networks exhibit the “small-world” phenomenon, i.e., any two vertices in the network is often connected by a small number of intermediate vertices (the shortest-path distance is small). The small-world phenomenon in the real populations was first discovered by Milgram [12]. In his study, the average distance between two Americans is around 6. Several recent studies [7, 10, 13] offer significant evidence to support similar observations in the online social networks and Internet itself. In addition, the power-law degree distribution (or scale-free property) which

many of these networks also directly lead to the small average distance [1]. Clearly, the small-world property can help facilitate the communication and speed up the diffusion process and information spreading in a large network.

However, the small-world effect can be a dangerous *double-edged sword*. When a system is benefited from the efficient communication and fast information diffusion, it also makes itself more vulnerable to various attacks: diseases, (computer) virus, spams, and misinformation, etc. For instance, it has been shown that a small-world graph can have much faster disease propagation than a regular lattice or a random graph [14]. Indeed, the six degrees of separation may suggest that a highly infectious disease could spread to all six billion people living in the earth about size incubation periods of the diseases [14]. The small-world property of Internet and WWW not only enables the computer virus and spams to be much easier to spread, but also makes them hard to stop. More recently, the misinformation problem in the social networks has made several public outcry [3]. These small-world online social network potentially facilitate the spread of misinformation to reach a large number of audience in short time, which may cause public panic and have other disruptive effects.

To prevent an outbreak, the most basic strategy is to remove the affected individuals (or computers) from the network system, like quarantine. However, in many situations, the explicit quarantine may be hard to achieve: the contagious individuals are either unknown or hard to detect; or it is often impossible to detect and remove each infected individual; or there are many already being affected and it become too costly to remove all of them in a timely fashion. Thus, it is important to consider alternative strategies to help prevent and even delay the spreading where the latter can be essential in discovering and/or deploying new methods for dealing with the outbreaks.

Recently, there have been a lot of interests in understanding the network factors (such as the small-world and scale-free properties) in the epidemics and information diffusion process, and utilizing the network structures in detecting/preventing the outbreaks. Several studies have focused on modeling the disease epidemics on the small-world and/or scale-free networks [14–16]; in [11], Leskovec *et al.* study how to deploy sensors cost-effectively in a network system (sensors are assigned to vertices) to detect an outbreak; in [3], Budak *et al.* consider how to limit the misinformation by identifying a set of individuals that are needed to adopt the “good” information (being immune in epidemics) in order to minimize those being affected by the “bad” information (being infected in epidemics). In addition, we note that from a different angle (viral marketing), there have been a list of studies on the *influence maximization* problem [8, 17], which aim to discover a set of most influential seeds to maximize the information spreading in the network. From the disease epidemics perspective, those seeds (assuming being selected using contagious model) may need particular protection to prevent an outbreak.

In this work, we study another basic and practically important strategy to help prevent and/or delay an outbreak in the context of network: limiting the contact between individuals. Different from the pure quarantine approach, here

individuals can still perform in the network system, though some contact relationships are forbidden. In other words, instead of removing vertices (individuals) form a network as in the quarantine approach, this strategy focuses on removing edges so that the (potential) outbreaks can be slowed down. Intuitively, if an individual contacts less number of other individuals, the chance for him or her to spread or being infected from the disease (misinformation) becomes less. From the network viewpoint, the edge-removal strategy essentially make the underlying (social) network less small-world, or simply “de-small-world”, i.e., the distances between individuals increase to delay the spreading process. In many situations, such a strategy is often easily and even voluntarily adopted. For instance, during the SARS epidemic in Beijing, 2004, there are much less people appearing in the public places. This approach can also be deployed in complement to the quarantine approach.

**Our Contribution.** Even though the edge-removal or *de-small-world* approach seems to be conceptually easy to understand, its mathematical foundation is still lack of study. Clearly, different edges (interactions) in the network are not being equivalent in terms of slowing down any potential outbreak: for a given individual, a link to an individual of high degree connection can be more dangerous than a link to another one with low degree connection. The edge importance (in terms of distance) especially coincides with Kleighnberg’s theoretical model [9] which utilizes the *long-range edges* on top of an underlying grid for explaining the small-world phenomenon. In this model, the long-range edges are the main factors which help connect the otherwise long-distance pairs with a smaller number of edges. However, there are no direct studies in fitting such a model to the real world graph to discover those long-range edges. In the mean time, additional constraint, such as the number of edges can be removed from the network, may exist because removing an edge can associate with certain cost. These factors and requirements give arise to the following fundamental research problem: *how can we maximally de-small-world a graph (making a graph to be less small-world) by removing a fixed number of edges?*

To tackle the problem, we make the following contributions:

1. We introduce the *average neighborhood size* as a new measure for the degree of being small-world and utilizes it to formally define the de-small-world network problem. Note that the typical average distance for measuring the small-world effects cannot uniformly treat the connected and disconnected networks; neither does it fit well with the spreading process. We also reformulate the de-small-world as the *local-reachable pair cut* problem.
2. We prove the NP-hardness of the general reachable pair cut problem and propose a greedy edge betweenness based approach as the benchmark in selecting the candidate edges for solving the de-small-world network. We transform the de-small-world network problem and express it as a OR-AND Boolean function maximization problem, which is also an NP-hard problem.

3. We develop a numerical relaxation approach to solve the de-small-world problem using its OR-AND boolean format. Our approach can find a local minimum based on the iterative gradient optimization procedure. In addition, we further generalize the betweenness measure and introduces the *short betweenness*, which measures the edge importance in terms of all the paths with distance no greater than a certain threshold. Using this measure, we can speed up the numerical relaxation approach by selecting a small set of candidate edges for removal.
4. We perform a detailed experimental evaluation, which demonstrates the effectiveness and efficiency of proposed approaches.

## 2 Problem Definition and Preliminary

In this section, we first formally define the *de-small-world* network problem and prove its NP-hardness; then we introduce the basic greedy approaches based on edge betweenness which will serve as the basic benchmark; and finally we show the de-small-world network problem can be transformed and expressed as a OR-AND Boolean function maximization problem.

**Problem Formulation.** In order to model the edge-removal process and formally define the *de-small-world* network problem, a criterion is needed to precisely capture the *degree* of being small-world. Note that here the goal is to help prevent and/or delay the potential outbreak and epidemic process. The typical measure of small-world network is based on the average distance (the average length of the shortest path between any pair of vertices in the entire network). However, this measure is not able to provide unified treatment of the connected and cut network. Specifically, assuming a connected network is broken into several cut network and the average distance on the cut network is not easy to express. On the other hand, we note that the de-small-world network graph problem is different from the network decomposition (clustering) problem which tries to break the entire network into several components (connected subgraphs). From the outbreak prevention and delaying perspective, the cost of network decomposition is too high and may not be effective. This is because each individual component itself may still be small-world; and the likelihood of completely separating the contagious/infected group from the rest of populations (the other components) is often impossible.

Given this, we introduce the *average neighborhood size* as a new measure for the degree of being small-world and utilize it to formally define the de-small-world network problem. Especially, the new measure can not only uniformly treat both connected and cut networks and aims to directly help model the spreading/diffusion process. Simply speaking, for each vertex  $v$  in a network  $G = (V, E)$  where  $V$  is the vertex set and  $E$  is the edge set, we define the neighborhood of  $v$  as the number of vertices with distance no greater than  $k$  to  $v$ , denoted as  $N^k(v)$ . Here  $k$  is the user-specified *spreading* (or delaying) parameter which aims to measure the outbreak speed, i.e., in a specified time unit, the maximum distance between individual  $u$  (source) to another one  $v$  (destination)

who can be infected if  $u$  is infected. Thus, the average neighborhood size of  $G$ ,  $\sum_{v \in V} N^k(v)$ , can be used to measure the robustness of the network with respect to a potential outbreak in a certain time framework. Clearly, a potential problem of the small-world network is that even for a small  $k$ , the average neighborhood size can be still rather large, indicating a large (expected) number of individuals can be quickly affected (within time framework  $k$ ) during an outbreak process.

Formally, the *de-small-world* network problem is defined as follows:

**Definition 1 (De-Small-World Network Problem).** Given the edge-removal budget  $L > 0$  and the spreading parameter  $k > 1$  we seek a subset of edges  $E_r \subset E$ , where  $|E_r| = L$ , such that the average neighborhood size is minimized:

$$\min_{|E_r|=L} \frac{\sum_{v \in V} N^k(v|G \setminus E_r)}{|V|}, \quad (1)$$

where  $N^k(v|G \setminus E_r)$  is the neighborhood size of  $v$  in the graph  $G$  after removing all edges in  $E_r$  from the edge set  $E$ .

Note that in the above definition, we assume each vertex has the equal probability to be the source of infection. In the general setting, we may consider to associate each vertex  $v$  with a probability to indicate its likelihood to be (initially) infected. Furthermore, we may assign each edge with a weight to indicate the cost to removing such an edge. For simplicity, we do not study those extensions in this work; though our approaches can be in general extended to handle those additional parameters. In addition, we note that in our problem, we require the spreading parameter  $k > 1$ . This is because for  $k = 1$ , this problem is trivial: the average neighborhood size is equivalent to the average vertex degree; and removing any edge has the same effect. In other words, when  $k = 1$ , the neighborhood criterion does not capture the spreading or cascading effects of the small-world graph. Therefore, we focus on  $k > 1$ , though in general  $k$  is still relatively small (for instance, no higher than 3 or 4 in general).

**Reachable Pair Cut Formulation:** We note the de-small-world network problem can be defined in terms of the *reachable pair cut* formulation. Let a pair of two vertices whose distance is no greater than  $k$  is referred to as a *local-reachable pair* or simply *reachable pair*. Let  $\mathcal{R}_G$  record the set of all local reachable pairs in  $G$ .

**Definition 2 (Reachable Pair Cut Problem).** For a given local  $(u, v)$ , if  $d(u, v|G) \leq k$  in  $G$ , but  $d(u, v|G \setminus E_s) > k$ , where  $E_s$  is an edge set in  $G$ , then we say  $(u, v)$  is being **local cut** (or simply cut) by  $E_s$ . Given the edge-removal budget  $L > 0$  and the spreading parameter  $k > 1$ , the reachable pair cut problem aims to find the edge set  $E_r \subseteq E$ , such that the maximum number of pairs in  $\mathcal{R}_G$  is cut by  $E_r$ .

Note that here the (local) cut for a pair of vertices simply refers to increase their distance; not necessarily completely disconnect them in the graph  $(G \setminus E_s)$ .

Also, since  $\mathcal{R}_{G \setminus E_r} \subseteq \mathcal{R}_G$ , i.e., every local-reachable pair in the remaining network  $G \setminus E_r$  is also the local-reachable in the original graph  $G$ , the problem is equivalently to maximize  $|\mathcal{R}_G| - |\mathcal{R}_{G \setminus E_r}|$  and minimize the number of local reachable pairs  $|\mathcal{R}_{G \setminus E_r}|$ . Finally, the correctness of such a reformulation (de-small-world problem = reachable pair cut problem) follows this simple observation:  $\sum_{v \in V} N^k(v|G) = 2|\mathcal{R}_G|$  (and  $\sum_{v \in V} N^k(v|G \setminus E_r) = 2|\mathcal{R}_{G \setminus E_r}|$ ). Basically, every reachable pair is counted twice in the neighborhood size criterion.

In the following, we study the hardness of the general reachable pair cut problem.

**Theorem 1.** *Given a set  $RS$  of local reachable pairs in  $G = (V, E)$  with respect to  $k$ , the problem of finding  $L$  edges  $E_r \subseteq E$  ( $|E_r| = L$ ) in  $G$  such that the maximal number of pairs in  $RS$  being cut by  $E_r$  is NP-Hard.*

All the proofs of Theorems and Lemmas can be found in our Appendix<sup>1</sup>. Note that in the general problem,  $RS$  can be any subset of  $\mathcal{R}_G$ . The NP-hardness of the general reachable pair cut problem a strong indicator that the de-small-world network problem is also hard. In addition, we note that the submodularity property plays an important role in solving vertex-centered maximal influence [8], outbreak detection [11], and limiting misinformation spreading [3] problems. However, such property does not hold for the edge-centered de-small-world problem.

**Lemma 1.** *Let set function  $f : 2^E \rightarrow Z^+$  records the number of local reachable pairs in  $\mathcal{R}_G$  is cut by an edge set  $E_s$  in graph  $G$ . Function  $f$  is neither submodular (diminishing return) nor supermodular.*

**Greedy Betweenness-Based Approach.** Finding the optimal solution for the de-small-world problem is likely to be NP-hard. Clearly, it is computationally prohibitive to enumerate all the possible removal edge set  $E_r$  and to measure how many reachable pairs could be cut or how much the average neighborhood size is reduced. In the following, we describe a greedy approach to heuristically discover a solution edge-set. This approach also serves as the benchmark for the de-small-world problem.

The basic approach is based on the edge-betweenness, which is a useful criterion to measure the edge importance in a network. Intuitively, the edge-betweenness measures the edge important with respect to the shortest paths in the network. The high betweenness suggests that the edge is involved into many shortest paths; and thus removing them will likely increase the distance of those pairs linked by these shortest paths. Here, we consider two variants of edge-betweenness: the (global) edge-betweenness [4] and the local edge-betweenness [5]. The global edge-betweenness is the original one [4] and is defined as follows:

$$B(e) = \sum_{s \neq t \in V} \frac{\delta_{st}(e)}{\delta_{st}},$$

<sup>1</sup> <https://www.dropbox.com/s/rpkqpn6y7mwconk/Appendix.pdf?dl=0>.

where  $\delta_{st}$  is the total number of shortest paths between vertex  $s$  and  $t$ , and  $\delta_{st}(e)$  the total number of shortest paths between  $u$  and  $v$  containing edge  $e$ .

The local edge-betweenness considers only those vertex pairs whose shortest paths are no greater than  $k$ , and is defined as

$$LB(e) = \sum_{s \neq t \in V, d(s,t) \leq k} \frac{\delta_{st}(e)}{\delta_{st}},$$

The reason to use the local edge-betweenness measure is because in the de-small-world (and reachable pair cut) problem, we focus on those local reachable pairs (distance no greater than  $k$ ). Thus, the contribution to the (global) betweenness from those pairs with distance greater than  $k$  can be omitted. The exact edge-betweenness can be computed in  $O(nm)$  worst case time complexity [2] where  $n = |V|$  (the number of vertices) and  $m = |E|$  (the number of edges) in a given graph, though in practical the local one can be computed much faster.

Using the edge-betweenness measure, we may consider the following *generic procedure* to select the  $L$  edges for  $E_r$ :

- (1) Select the top  $r < L$  edges into  $E_r$ , and remove those edges from the input graph  $G$ ;
- (2) Recompute the betweenness for all remaining edges in the updated graph  $G$ ;
- (3) Repeat the above procedure  $\lceil L/r \rceil$  times until all  $L$  edges are selected.

Note that the special case  $r = 1$ , where we select each edge in each iteration, the procedure is very similar to the Girvan-Newman algorithm [4] in which they utilize the edge-betweenness for community discovery. Gregory [5] generalizes it to use the local-edge betweenness. Here, we only consider to pickup  $L$  edges and allow users to select the frequency to recompute the edge-betweenness (mainly for efficiency consideration). The overall time complexity of the betweenness based approach is  $O(\lceil L/r \rceil nm)$  (assuming the exact betweenness computation is adopted).

## 2.1 OR-AND Boolean Function and Its Maximization Problem

In the following, we transform the de-small-world network problem and express it as a OR-AND Boolean function maximization problem, which forms the basis for our optimization problem in next section. First, we will utilize the OR-AND graph to help represent the de-small-world (reachable pair cut) problem. Let us denote  $P$  the set of all the short paths in  $G$  that have length at most  $k$ .

**OR-AND Graph:** Given a graph  $G = (V, E)$ , the vertex set of an *OR-AND* graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  is comprised of three kinds of nodes  $\mathcal{V}_E$ ,  $\mathcal{V}_P$  and  $\mathcal{V}_{\mathcal{R}_G}$ , where each node in  $\mathcal{V}_E$  corresponds to a unique edge in  $E$ , each node in  $\mathcal{V}_P$  corresponds to a short path in  $P$ , and each node in  $\mathcal{V}_{\mathcal{R}_G}$  corresponds to a unique reachable pair in  $G$  (with respect to  $k$ ). Figure 1(b) shows those nodes for graph  $G$  in Fig. 1(a). The edge set consists of two types of edges: (1) Each short path node

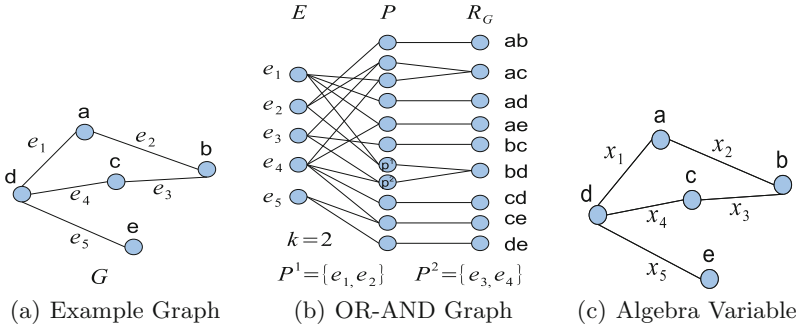


Fig. 1. OR-AND graph and algebra variable

in  $\mathcal{V}_P$  is linked with the vertices in  $\mathcal{V}_E$  corresponding to those edges in the path. For instance path node  $p^1$  in  $\mathcal{V}_P$  links to edge node  $e_1$  and  $e_2$  in  $\mathcal{V}_E$  in Fig. 1(b). Each reachable pair node in  $\mathcal{V}_{R_G}$  links to those path nodes which connects the reachable pair. For instance, the reachable pair  $bd$  is connected with path node  $p^1$  and  $p^2$  in Fig. 1(b).

Intuitively, in the OR-AND graph, we can see that in order to cut a reachable pair, we have to cut *all* the short paths between them (AND). To cut one short path, we need to remove only one edge in that path (OR). Let  $P(u, v)$  consists all the (simple) short paths between  $u$  and  $v$  whose length are no more than  $k$ . For each short path  $p$  in  $P(u, v)$ , let  $e$  corresponds to a Boolean variable for edge  $e \in p$ : if  $e_i = T$ , then the edge  $e_i$  is not cut; if  $e_i = F$ , then the edge is cut ( $e_i \in E_r$ ). Thus, for each reachable pair  $(u, v) \in \mathcal{R}_G$ , we can utilize the a Boolean OR-AND expression to describe it:

$$I(u, v) = \bigvee_{p \in P(u, v)} \bigwedge_{e \in p} e \tag{2}$$

For instance, in the graph  $G$  (Fig. 1(b)),

$$I(b, d) = (e_1 \wedge e_2) \vee (e_3 \wedge e_4)$$

Here,  $I(b, d) = T$  indicating the pair is being cut only if for both  $p^1$  and  $p^2$  are cut. For instance, if  $e_1 = F$  and  $e_3 = F$ , then  $I(b, d) = F$ ; and  $e_1 = F$ , but  $e_3 = T$  and  $e_4 = T$ ,  $I(b, d) = T$ . Given this, the de-small-world problem (and the reachable pair cut problem) can be expressed as the following Boolean function maximization problem.

**Definition 3 (Boolean Function Maximization Problem).** Given a list of Boolean functions (such as  $I(u, v)$ , where  $(u, v) \in \mathcal{R}_G$ ), we seek a Boolean variable assignment where exactly  $L$  variables are assigned false ( $e = F$  iff  $e \in E_r$ , and  $|E_r| = L$ ), such that the maximal number of Boolean functions being false ( $I(u, v) = F$  corresponding to  $(u, v)$  is cut by  $E_r$ ).



Unfortunately, the Boolean function maximization problem is also NP-hard since it can directly express the general reachable pair cut problem. In the next section, we will introduce a numerical relation approach to solve this problem.

### 3 Path Algebra and Optimization Algorithm

In this section, we introduce a numerical relaxation approach to solve the Boolean function maximization problem (and thus the de-small-world problem). Here, the basic idea is that since the direct solution for the Boolean function maximization problem is hard, instead of working on the Boolean (binary) edge variable, we relax to it to be a numerical value. However, the challenge is that we need to define the numerical function optimization problem such that it meet the following two criteria: (1) it is rather accurately match the Boolean function maximization; and (2) it can enable numerical solvers to be applied to optimize the numerical function. In Subsect. 3.1, we introduce the numerical optimization problem based on the path algebra. In Subsect. 3.2, we discuss the optimization approach for solving this problem.

#### 3.1 Path-Algebra and Numerical Optimization Problem

To construct a numerical optimization problem for the Boolean function maximization format of the de-small-world problem, we introduce the following path-algebra to describe all the short paths between any reachable pair in  $\mathcal{R}_G$ . For each edge  $e$  in the graph  $G = (V, E)$ , we associate it with a variable  $x_e$ . Then, for any reachable pair  $(u, v) \in \mathcal{R}_G$ , we define its corresponding path-algebra expression  $\mathcal{P}(u, v)$  as follows:

$$\mathcal{P}(u, v) = \sum_{p \in P(u, v)} \prod_{e \in p} x_e \quad (3)$$

Taking the path-algebra for  $(b, d)$  in Fig. 1 and (c) as example, we have

$$\mathcal{P}(b, d) = x_2x_1 + x_3x_4$$

Basically, the path-algebra expression  $\mathcal{P}(u, v)$  directly corresponds to the Boolean expression  $I(u, v)$  by replacing *AND*( $\wedge$ ) with product ( $\times$ ), *OR*( $\vee$ ) with sum ( $+$ ), and Boolean variable  $e$  with algebraic variable  $x_e$ . Intuitively,  $\mathcal{P}(u, v)$  records the weighted sum of each path in  $P(u, v)$ , where the weight is the product based on the edge variable  $x_e$ . Note that when  $x_e = 1$  for every edge  $e$ , when  $\mathcal{P}(u, v)$  simply records the number of different short paths (with length no more than  $k$ ) between  $u$  and  $v$ , i.e.,  $\mathcal{P}(u, v) = |P(u, v)|$ . Furthermore, if assuming  $x_e \geq 0$ , then  $\mathcal{P}(u, v) = 0$  is equivalent to in each path  $p \in P(u, v)$ , there is at least one edge variable is equivalent to 0. In other words, assuming if variable  $x_e = 0$  iff  $e = T$ , then  $\mathcal{P}(u, v) = 0$  iff  $I(u, v) = F$  and  $\mathcal{P}(u, v) > 0$  iff  $I(u, v) = T$ .

Given this, we may be tempted to optimize the follow objective function based on the path-algebra expression to represent the Boolean function maximization problem:

$\sum_{(u,v) \in \mathcal{R}_G} \mathcal{P}(u, v)$ . However, this does not accurately reflect our goal, as to minimize  $\sum_{(u,v) \in \mathcal{R}_G} \mathcal{P}(u, v)$ , we may not need any  $\mathcal{P}(u, v) = 0$  (which shall be our main goal). This is because  $\mathcal{P}(u, v)$  corresponds to the weighted sum of path products. Can we use the path-algebra to address the importance of  $\mathcal{P}(u, v) = 0$  in the objective function?

We provide a positive answer to this problem by utilizing an exponential function transformation. Specifically, we introduce the following *numerical maximization* problem based on the path expression:

$$\sum_{(u,v) \in \mathcal{R}_G} e^{-\lambda \mathcal{P}(u,v)}, \text{ where, } 0 \leq x_e \leq 1, \sum x_e \geq X - L \tag{4}$$

Note that  $0 \leq e^{-\lambda \mathcal{P}(u,v)} \leq 1$  (each  $x_e \geq 0$ ), and only when  $\mathcal{P}(u, v) = 0$ ,  $e^{-\lambda \mathcal{P}(u,v)} = 1$  (the largest value for each term). When  $\mathcal{P}(u, v) \approx 1$ , the term  $e^{-\lambda \mathcal{P}(u,v)}$  can be rather small (approach 0). The parameter  $\lambda$  is the adjusting parameter to help control the exponential curve and smooth the objective function. Furthermore, the summation constraint  $\sum x_e \geq X - L$  is to express the budget condition that there shall have  $L$  variables with  $x_i \approx 0$ . Here  $X$  is the total number of variables in the objective function ( $X = |E|$  if we consider every single edge variable  $x_e$ ).

### 3.2 Gradient Optimization

Clearly, it is very hard to find the exact (or closed form) solution for maximizing function in Eq. 4 under these linear constraints. In this section, we utilize the standard *gradient* (ascent) approach together with the *active set* method [6] to discover a local maximum. The gradient ascent takes steps proportional to the positive of the gradient iteratively to approach a local minimum. The active set approach is a standard approach in optimization which deals with the *feasible regions* (represented as constraints). Here we utilize it to handle the constraint in Eq. 4.

**Gradient Computation:** To perform gradient ascent optimization, we need compute the gradient  $g(x_e)$  for each variable  $x_e$ . Fortunately, we can derive a closed form of  $g(x_e)$  in  $\sum_{(u,v) \in \mathcal{R}_G} e^{-\lambda \mathcal{P}(u,v)}$  as follows:

$$g(x_e) = \frac{\partial \sum_{(u,v) \in \mathcal{R}_G} e^{-\lambda \mathcal{P}(u,v)}}{\partial x_e} = \sum_{(u,v) \in \mathcal{R}_G} -\lambda \mathcal{P}(u, v, e) e^{-\lambda \mathcal{P}(u,v)},$$

where  $\mathcal{P}(u, v, e)$  is the sum of the path-product on all the paths going through  $e$  and we treat  $x_e = 1$  in the path-product. More precisely, let  $P(u, v, e)$  be the set of all short paths (with length no more than  $k$ ) between  $u$  and  $v$  going through edge  $e$ , and then,

$$\mathcal{P}(u, v, e) = \sum_{p \in P(u,v,e)} \prod_{e' \in p \setminus \{e\}} x_{e'} \tag{5}$$

Using the example in Fig. 1 and (c), we have

$$\mathcal{P}(b, d, e_1) = x_2$$

Note that once we have all the gradients for each edge variable  $x_e$ , then we update them accordingly,

$$x_e = x_e + \beta g(x_e),$$

where  $\beta$  is the step size (a very small positive real value) to control the rate of convergence.

**$\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$  Computation** To compute the gradient, we need compute all  $\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$  for  $(u, v) \in \mathcal{R}_G$ . Especially, the difficulty is that even compute the total number of simple short paths (with length no more than  $k$ ) between  $u$  and  $v$ , denoted as  $|P(u, v)|$  is known to be expensive. In the following, we describe an efficient procedure to compute  $\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$  efficiently. The basic idea is that we perform a DFS from each vertex  $u$  with traversal depth no more than  $k$ . During the traversal from vertex  $u$ , we maintain the partial sum of both  $\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$  for each  $v$  and  $e$  where  $u$  can reach within  $k$  steps. After each traversal, we can then compute the exact value of  $\mathcal{P}(u, *)$  and  $\mathcal{P}(u, *, *)$ .

The DFS procedure starting from  $u$  to compute all  $\mathcal{P}(u, *)$  and  $\mathcal{P}(u, *, *)$  is illustrated in Algorithm 1 (Appendix D). In Algorithm 1, we maintain the current path (based on the DFS traversal procedure) in  $p$  and its corresponding product  $\sum_{e \in p} x_e$  is maintained in variable  $w$  (Line 9 and 10). Then, we incrementally update  $\mathcal{P}(u, v)$  assuming  $v$  is the end of the path  $p$  (Line 11). In addition, we go over each edge in the current path, and incrementally update  $\mathcal{P}(u, v)$  ( $w/x_e = \prod_{e' \in p \setminus \{e\}} x_{e'}$ , Line 13.) Note that we need invoke this procedure for every vertex  $u$  to compute all  $\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$ . Thus, the overall time complexity can be written as  $O(|V|\bar{d}^k)$  for a random graph where  $\bar{d}$  is the average vertex degree.

The overall gradient optimization algorithm is depicted in Algorithm 2 (Appendix E). Here, we use  $\mathcal{C}$  to describe all the edges which need be processed for optimization. At this point, we consider all the edges and thus  $\mathcal{C} = E$ . Later, we will consider to first select some candidate edges. The entire algorithm performs iteratively and each iteration has three major steps:

**Step 1** (Lines 6–8): it calculates the gradient  $g(x_e)$  of for every edge variable  $x_e$  and an average gradient  $\bar{g}$ ;

**Step 2** (Lines 9–16): only those variables are not in the active set  $\mathcal{A}$  will be updated. Specifically, if the condition ( $\sum_{e \in E} x_e \geq |E| - L$ ) is not met, we try to adjust  $x_e$  back to the feasible region. Note that by using  $g(x_e) - \bar{g}$  (Line 11) instead of  $g(x_e)$  (Line 13), we are able to increase the value of those  $x_e$  whose gradient is below average. However, such adjustment can still guarantee the overall objective function is not decreased (thus will converge). Also, we make sure  $x_e$  will be between 0 and 1.

**Step 3** (Lines 17–22): the active set is updated. When an edge variable reaches 0 or 1, we put them in the active set so that we will not need to update them in Step 2. However, for those edges variables in the active set, if their gradients are less (higher) than the average gradient for  $x_e = 0$  ( $x_e = 1$ ), we will release them from the active set and let them to be further updated.

Note that the gradient ascent with the active set method guarantees the convergence of the algorithm (mainly because the overall objective function is not decreased). However, we note that in Algorithm 2, the bounded condition ( $\sum_{e \in E} x_e \geq |E| - L$ ) may not be necessarily satisfied even with the update in Line 11. Though this can be achieved through additional adjustment, we do not consider them mainly due to the goal here is not to find the exact optimization, but mainly on identifying the smallest  $L$  edges based on  $x_e$ . Finally, the overall time complexity of the optimization algorithm is  $O(t(|V| * \bar{d}^k + |E|))$ , given  $t$  is the maximum number of iterations before convergence.

## 4 Short Betweenness and Speedup Techniques

In Sect. 3, we reformulate our problem into a numerical optimization problem. We further develop an iterative *gradient* algorithm to select the top  $L$  edges in to  $E_r$ . However, the basic algorithm can not scale well to very large graphs due to the large number ( $|E|$ ) of variables involved. In this section, we introduce a new variant of the edge-betweenness and use it to quickly reduce the variables needed in the optimization algorithm (Algorithm 2). In addition, we can further speedup the DFS procedure to compute  $\mathcal{P}(u, v)$  and  $\mathcal{P}(u, v, e)$  in Algorithm 1.

**Short Betweenness.** In this subsection, we consider the following question: *What edge importance measure can directly correlate with  $x_e$  in the objective function in Eq. 4 so that we can use it to help quickly identify a candidate edge set for the numerical optimization described in Algorithm 2?* In this work, we propose a new edge-betweenness measure, referred to as the *short betweenness* to address the this question. It is intuitively simple and has an interestingly relationship with respect to the gradient  $g(x_e)$  for each edge variable. It can even be directly applied for selecting  $E_r$  using the generic procedure in Sect. 2 and is much more effective compared with the global and local edge-betweenness which measure the edge importance in terms of the shortest path (See comparison in Sect. 5).

Here we formally define  $\nabla(e_i)$  as *short betweenness*.

**Definition 4 (Short Betweenness).** The short betweenness  $SB(e)$  for edge  $e$  is as follows,  $SB(e) = \sum_{(u,v) \in \mathcal{R}_G} \frac{|P(u,v,e)|}{|P(u,v)|}$ .

Recall that  $(u, v) \in \mathcal{R}_G$  means  $d(u, v) \leq k$ ;  $|P(u, v)|$  is the number of short paths between  $u$  and  $v$ ; and  $|P(u, v, e)|$  is the number of short paths between  $u$  and  $v$  which must go through edge  $e$ . The following lemma highlights the relationship between the short betweenness and the gradient of edge variable  $x_e$ :

**Lemma 2.** *Assuming for all edge variables  $x_e = 1$ , then  $g(x_e) \geq -SB(e)$ .*

Basically, when  $x_e = 1$  for every edge variable  $x_e$  (this is also the initialization of Algorithm 2), the (negative) short betweenness serves a lower bound of the gradient  $g(e)$ . Especially, since the gradient is negative, the higher the gradient of  $|g(e)|$  is, the more likely it can maximize the objective function (cut more reachable pairs in  $\mathcal{R}_G$ . Here, the short betweenness  $SB(e)$  thus provide an upper bound (or approximation) on  $|g(e)|$  (assuming all other edges are presented in the graph); and measures the the edge potential in removing those local reachable pairs. Finally, we note that Algorithm 1 can be utilized to compute  $|P(u, v)|$  and  $|P(u, v, e)|$ , and thus the short betweenness (just assuming  $x_e = 1$  for all edge variables).

**Scaling Optimization Using Short Betweenness:** First, we can directly utilize the short betweenness to help us pickup a candidate set of edge variables, and then Algorithm 2 only need to work on these edge variables (considering other edge variables are set as 1). Basically, we can choose a subset of edges  $E_s$  which has the highest short betweenness in the entire graph. The size of  $E_s$  has to be larger than  $L$ ; in general, we can assume  $|E_s| = \alpha L$ , where  $\alpha > 1$ . In the experimental evaluation (Sect. 5), we found when  $\alpha = 5$ , the performance of using candidate set is almost as good as the original algorithm which uses the entire edge variables. Once the candidate set edge set is selected, we make the following simple observation:

**Lemma 3.** *Given a candidate edge set  $E_s \subseteq E$ , if any reachable pair  $(u, v) \in \mathcal{R}_G$  can be cut by  $E_r$  where  $E_r \subseteq E_s$  and  $|E_r| = L$ , then, each path in  $P(u, v)$  must contains at least one edge in  $E_s$ .*

Clearly, if there is one path in  $P(u, v)$  does not contain an edge in  $E_s$ , it will always linked no matter how we select  $E_r$  and thus cannot cut by  $E_r \subseteq E_s$ . In other words,  $(u, v)$  has to be cut by  $E_s$  if it can be cut by  $E_r$ . Given this, we introduce  $\mathcal{R}_s = \mathcal{R}_G \subseteq \mathcal{R}_{G \setminus E_s}$ . Note that  $\mathcal{R}_s$  can be easily computed by the DFS traversal procedure similar to Algorithm 1. Thus, we can focus on optimizing

$$\sum_{(u,v) \in \mathcal{R}_s} e^{-\lambda P(u,v)}, \text{ where, } 0 \leq x_e \leq 1, \sum x_e \geq X - L \quad (6)$$

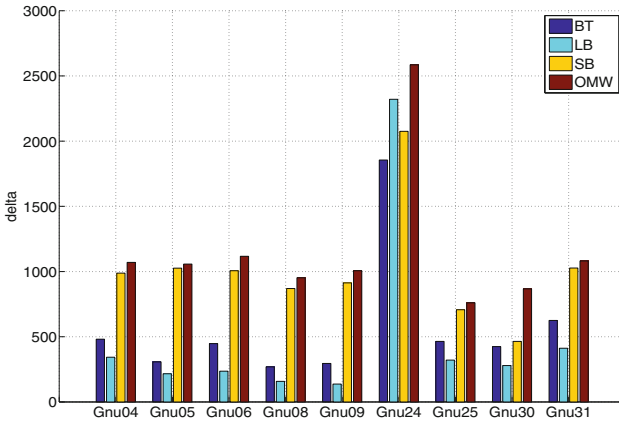
Furthermore, let  $E_P = \bigcup_{(u,v) \in \mathcal{R}_s} \bigcup_{p \in P(u,v)} p$ , which records those edges appearing in certain path linking a reachable pair cut by  $E_P$ . Clearly, for those edges in  $E \setminus E_P$ , we can simply prune them from the original graph  $G$  without affecting the final results. To sum, the short betweenness measure can help speed up the numerical optimization process by reducing the number of edge variables and pruning non-essential edges from the original graph.

## 5 Experimental Study

In this section, we report the results of the empirical study of our methods. Specifically, we are interested in the performance (in terms of reachable pair cut) and the efficiency (running time).

**Table 1.** Network statistics

| Dataset    | $ V $  | $ E $   | $\pi$ |
|------------|--------|---------|-------|
| Gnutella04 | 10,876 | 39,994  | 9     |
| Gnutella05 | 8,846  | 31,839  | 9     |
| Gnutella06 | 8,717  | 31,525  | 9     |
| Gnutella08 | 6,301  | 20,777  | 9     |
| Gnutella09 | 8,114  | 26,013  | 9     |
| Gnutella24 | 26,518 | 65,369  | 10    |
| Gnutella25 | 22,687 | 54,705  | 11    |
| Gnutella30 | 36,682 | 88,328  | 10    |
| Gnutella31 | 62,586 | 147,892 | 11    |



**Fig. 2.**  $\delta$  for all real datasets

**Performance:** Given a set of edges  $E_r$  with budget  $L$ , the total number of reachable pairs being cut by  $E_r$  is  $|\mathcal{R}_G| - |\mathcal{R}_{G \setminus E_r}|$  or simply  $\Delta|\mathcal{R}_G|$ . We use the average pair being cut by an edge, i.e.,  $\delta = \frac{\Delta|\mathcal{R}_G|}{L}$  as the performance measure.

**Efficiency:** The running time of different algorithms.

**Methods:** Here we compare the following methods:

- (1) *Betweenness based method*, which is defined in terms of the shortest paths between any two vertices in the whole graph  $G$ ; hereafter, we use  $BT$  to denote the method based on this criterion.
- (2) *Local Betweenness based method*, which, compared with betweenness method( $BT$ ), takes only the vertex pair within certain distance into consideration; hereafter, we use  $LB$  to stand for the method based on local betweenness.

- (3) *Short Betweenness based method*, the new betweenness introduced in this paper, which considers all short paths whose length is no more than certain threshold. Here we denote the method based on short betweenness as *SB*.
- (4) *Numerical Optimization method*, which solves the de-small-world problem iteratively by calculating gradients and updating the edge variables  $x_e$ . Based on whether the method use the candidate set or not, we have two versions of optimization methods: *OMW* (Optimization Method With candidate set) and *OMO* (Optimization Method withOut candidate set). Note that we normally choose the top  $5L$  edges as our candidate set.

We have a generic procedure to select  $L$  edges depending on parameter  $r$  (batch size) (Sect. 2). We found for different methods *BT*, *LB* and *SB*, the effects of  $r$  seem to be rather small (as illustrated in Fig. 3). Thus, in the reminder of the experiments, we choose  $r = L$ , i.e., we select the top  $L$  edges using the betweenness calculated for the entire (original graph).

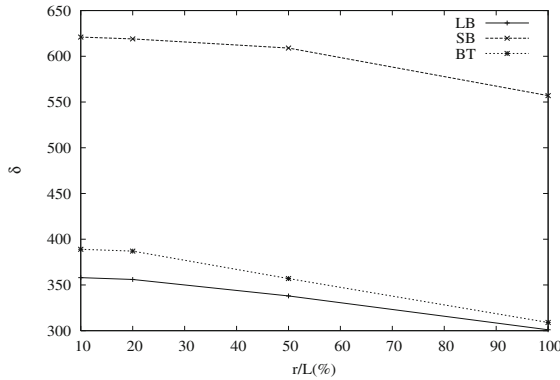


Fig. 3. Varying  $\frac{r}{L}$

Table 2. Time (Seconds)

| Time   | <i>BT</i> | <i>LB</i> | <i>SB</i> | OMW     |
|--------|-----------|-----------|-----------|---------|
| 10,876 | 382.27    | 24.82     | 33.75     | 1021.66 |
| 8,846  | 21346.54  | 496.17    | 8.98      | 110.80  |
| 62586  | 392.54    | 25.31     | 34.60     | 1092.55 |

Table 3.  $\delta$  By Varying  $l$

| $l$  | <i>BT</i> | <i>LB</i> | <i>SB</i> | OMW | OMO |
|------|-----------|-----------|-----------|-----|-----|
| 500  | 240       | 415       | 912       | 996 | 973 |
| 1000 | 261       | 372       | 740       | 803 | 805 |
| 2000 | 301       | 329       | 572       | 620 | 622 |

Table 4.  $\delta$  By Varying  $k$

| $k$ | <i>BT</i> | <i>LB</i> | <i>SB</i> | OMW  | OMO |
|-----|-----------|-----------|-----------|------|-----|
| 2   | 25        | 32        | 55        | 58   | 58  |
| 3   | 261       | 372       | 740       | 803  | 805 |
| 4   | 761       | 976       | 2113      | 2389 | -   |

All the algorithms are implemented using C++ and the Standard Template Library (STL), and the experiments are conducted on a 2.0 GHz Dual Core AMD Opteron CPU with 4.0 GB RAM running on Linux.

We study the performance of our algorithms on real datasets. The benchmarking datasets are listed in Table 1. All networks contain certain properties

commonly observed in social networks, such as small diameter. All datasets are downloadable from Stanford Large Network Dataset Collection<sup>2</sup>.

In Table 1, we present important characteristics of all real datasets, where  $\pi$  is graph diameter. All these nine networks are snapshots of the Gnutella peer to peer file sharing network starting from August 2002. Nodes stand for the hosts in the Gnutella network topology and the edges for the connections between the hosts (Table 2).

**Varying  $L$ :** We perform this group of experiments on dataset *Gnu05* and we fix  $k = 3$ . Here we run these methods on three different edge budget  $L$ : 500, 1000 and 2000. The result is reported in Table 3. The general trend is that with smaller  $L$ ,  $\delta$  becomes bigger. This is because the set of reachable pairs removed by different edges could have intersection; when one edge is removed, the set of reachable pairs for other edges is also reduced. For particular methods, *BT* and *OMO* methods produces the lowest and highest  $\delta$ , and the different between *OMW* and *OMO* is very small.

**Varying  $k$ :** In this group of experiments, we fix  $L = 1000$  and we choose *Gnu04*. Here we choose three values for  $k$ : 2, 3 and 4. The result is reported in Table 4. From the result, we can see that when  $k$  becomes bigger,  $\delta$  become higher. This is also reasonable: when  $k$  becomes bigger, more reachable pairs are generated and meanwhile  $|E|$  is constant; therefore, each edge is potentially able to remove more reachable pairs. From the above three groups of experiments, we can see that *OMO* does not produce significant results compared with *OMW*. Therefore, in the following experiment, we do not study *OMO* method again.

**$\delta$  on all Real Datasets:** In this groups of experiment, we study the performance of each method on these nine datasets, with  $L$  being proportional to  $|E|$ . Specifically,  $L = |E| \times 1\%$ . We report the result in Fig. 2. *LB* generally produces the lowest  $\delta$ , around half that of *BT*; and also the best method, is the *SB* and *OMW* methods. Specifically, *OMW* is always slightly better than *SB*.

## 6 Conclusion

In this paper, we introduce the *de-small-world* network problem; to solve it, we first present a greedy edge betweenness based approach as the benchmark and then provide a numerical relaxation approach to solve our problem using OR-AND boolean format, which can find a local minimum. In addition, we introduce the *short-betweenness* to speed up our algorithm. The empirical study demonstrates the efficiency and effectiveness of our approaches. In the future, we plan to utilize MapReduce framework (e.g. Hadoop) to scale our methods to handle graphs with tens of millions of vertices.

<sup>2</sup> <http://snap.stanford.edu/data/index.html>.



## References

1. Andersen, R., Chung, F., Lu, L.: Modeling the small-world phenomenon with local network flow. *Internet Math.* **2**, 359–385 (2005)
2. Brandes, U.: A faster algorithm for betweenness centrality. *J. Math. Sociol.* **25**, 163–177 (2001)
3. Budak, C., Agrawal, D., El Abbadi, A.: Limiting the spread of misinformation in social networks. In: *Proceedings of the 20th International Conference on World Wide Web, WWW 2011* (2011)
4. Girvan, M., Newman, M.E.J.: Community structure in social and biological networks. *Proc. Nat. Acad. Sci.* **99**(12), 7821–7826 (2002)
5. Gregory, S.: Local betweenness for finding communities in networks. Technical report, University of Bristol, February 2008
6. Hager, W.W., Zhang, H.: A new active set algorithm for box constrained optimization. *SIAM J. Optim.* **17**, 526–557 (2006)
7. Jin, S., Bestavros, A.: Small-world characteristics of internet topologies and implications on multicast scaling. *Comput. Netw.* **50**, 648–666 (2006)
8. Kempe, D., Kleinberg, J., Tardos, E.: Maximizing the spread of influence through a social network. In: *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2003*, pp. 137–146 (2003)
9. Kleinberg, J.: The small-world phenomenon: an algorithmic perspective. In: *32nd ACM Symposium on Theory of Computing*, pp. 163–170 (2000)
10. Leskovec, J., Horvitz, E.: Planetary-scale views on a large instant-messaging network. In: *Proceedings of the 17th International Conference on World Wide Web, WWW 2008* (2008)
11. Leskovec, J., Krause, A., Guestrin, C., Faloutsos, C., VanBriesen, J., Glance, N.: Cost-effective outbreak detection in networks. In: *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2007*, pp. 420–429 (2007)
12. Milgram, S.: The small world problem. *Psychol. Today* **2**, 60–67 (1967)
13. Mislove, A., Marcon, M., Gummadi, K.P., Druschel, P., Bhattacharjee, B.: Measurement and analysis of online social networks. In: *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement, IMC 2007* (2007)
14. Moore, C., Newman, M.E.J.: Epidemics and percolation in small-world networks. *Phys. Rev. E* **61**, 5678–5682 (2000)
15. Newman, M.E.J.: Spread of epidemic disease on networks. *Phys. Rev. E* **66**, 016128+ (2002)
16. Pastor-Satorras, R., Vespignani, A.: Epidemic spreading in scale-free networks. *Phys. Rev. Lett.* **86**, 3200–3203 (2001)
17. Richardson, M., Domingos, P.: Mining knowledge-sharing sites for viral marketing. In: *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2002*, pp. 61–70 (2002)