

Controlling Infection by Blocking Nodes and Links Simultaneously*

Jing He¹, Hongyu Liang¹, and Hao Yuan²

¹ Institute for Interdisciplinary Information Sciences, Tsinghua University
{he-j08, lianghy08}@mails.tsinghua.edu.cn

² Department of Computer Science, City University of Hong Kong
haoyuan@cityu.edu.hk

Abstract. In this paper we study the problem of controlling the spread of undesirable things (viruses, epidemics, rumors, etc.) in a network. We present a model called the *mixed generalized network security model*, denoted by $\text{MGNS}(d)$, which unifies and generalizes several well-studied infection control model in the literature. Intuitively speaking, our goal under this model is to secure a subset of nodes and links in a network so as to minimize the expected total loss caused by a possible infection (with a spreading limit of d -hops) plus the cost spent on the preventive actions. Our model has wide applications since it incorporates both node-deletion and edge-removal operations. Our main results are as follows:

1. For all $1 \leq d < \infty$, we present a polynomial time $(d+1)$ -approximation algorithm for computing the optimal solution of $\text{MGNS}(d)$. This improves the approximation factor of $2d$ obtained in [19] for a special case of our model. We derive an $O(\log n)$ -approximation for the case $d = \infty$. Moreover, we give a polynomial time $\frac{3}{2}$ -approximation for $\text{MGNS}(1)$ on bipartite graphs.
2. We prove that for all $d \in \mathbb{N} \cup \{\infty\}$, it is \mathcal{APX} -hard to compute the optimum cost of $\text{MGNS}(d)$ even on 3-regular graphs. We also show that, assuming the Unique Games Conjecture [13], we cannot obtain a $(\frac{3}{2} - \epsilon)$ -approximation for $\text{MGNS}(d)$ in polynomial time. Our hardness results hold for the special case $\text{GNS}(d)$ in [19] as well.
3. We show that an optimal solution of $\text{MGNS}(d)$ can be found in polynomial time for every fixed $d \in \mathbb{N} \cup \{\infty\}$ if the underlying graph is a tree, and the infection cost and attack probability are both uniform. Our algorithm also works for the case where there are budget constraints on the number of secured nodes and edges in a solution. This in particular settles an open question from [21] that asks whether there exists an efficient algorithm for the minimum average contamination problem on trees.

* The first two authors were supported in part by the National Basic Research Program of China Grant 2007CB807900, 2007CB807901, and the National Natural Science Foundation of China Grant 61033001, 61061130540, 61073174. Portions of this work were supported by a grant from City University of Hong Kong (Project No. 7200218).

1 Introduction

During the recent years, much effort has been devoted to the study on the structure of various types of networks such as social networks, wireless sensor networks, computer networks, transportation networks, and the World Wide Web. An important and active subject is to study the *information diffusion process* in the situations where we want some news, topics, thoughts or products to spread quickly in the network, such as viral marketing [8]. This idea is formalized by Kempe, Kleinberg and Tardos [12] as a combinatorial problem called the *influence maximization problem*, which has since then been extensively studied under various settings (see, e.g., [6,10,15,20]).

In contrast, another important line of research is to study how to prevent or limit the spread of undesirable things through the network, such as the propagation of computer viruses and worms over computer networks, the fast spreading of malicious rumors through social networks, and the spread of infections or epidemics (such as Swine Flu and H1N1) among groups of people. In all these circumstances we need to eliminate or at least control the evolution of the bad things over the whole network, which is usually achieved by taking some preventive measures before the emergence of these undesirable things, and isolating or restricting the behaviors of some individuals if the infection has already been spread through the network. An important issue in real-world applications is the balance between the cost spent on prevention and the expected loss caused by infection. For example, installing anti-virus softwares on the computers is a natural response to the possible virus attack, but it may cost a lot of money and bring inefficiency to the protected computers due to high maintenance cost or memory requirement.

An elegant model that integrates both the security and infection costs has been formalized by Aspnes, Chang and Yampolskiy [3]. In their model, we seek for a subset of nodes on which we shall install the anti-virus softwares (call such nodes *secure*). A virus-attack is initiated by choosing one node from the network uniformly at random, and this node, if not secure, will infect all other nodes that are reachable from it in the network with all secure nodes removed. The goal is to minimize the cost for installing softwares (*security cost*) plus the expected total loss caused by the virus (*infection cost*). They consider both centralized (optimization) and game-theoretic settings. The model is substantially generalized by Kumar et al. [19] by allowing individual security and infection costs and arbitrary distribution of the virus-attack probability, and by introducing a parameter d into the model that represents the distance within the network that an infection can spread. This new model is called the *generalized network security model*, denoted $GNS(d)$. Thus, $GNS(d)$ is able to capture networks with less infection power or limited local information, such as ad hoc wireless networks. An issue with $GNS(d)$ is that it lacks the power of modeling the action of restricting the interconnections between individuals in the network (instead of simply removing them from the network), which, in the graph language, corresponds to blocking edges in the graph instead of deleting nodes. In spirit of such consideration, the contamination minimization model where edges are supposed to be blocked is raised by [16] and has been further studied in, e.g., [17,18,21].

In this paper, we present a model for minimizing the spread of infection that unifies and further generalizes the two aforementioned approaches, which we call the *mixed generalized network security model*, denoted by $\text{MGNS}(d)$. In our model, each node has its own *security cost* and *infection cost* as in $\text{GNS}(d)$, and each edge has its own *link-blocking cost* that represents the lost caused by the removal of the edge. The attack probability distribution can be arbitrary as in $\text{GNS}(d)$. The insecure node that is attacked initially will infect exactly those nodes that are within distance at most d from it in the *attack graph* obtained by removing all secure nodes and blocked edges from the original network. The cost of a solution is equal to the total expected infection cost of the nodes plus the cost for securing nodes and blocking edges in this solution. The goal is then to find a solution with minimum cost. Our main results in this paper, some of which improve on the previously best known results achieved for special cases of our model, are given in the following.

1. For all $1 \leq d < \infty$, we present a polynomial time $(d + 1)$ -approximation algorithm for computing the optimal solution of $\text{MGNS}(d)$ based on the primal-dual method. This improves the approximation factor of $2d$ obtained in [19] for $\text{GNS}(d)$, which is a special case of $\text{MGNS}(d)$. (We note that it is possible to design a reduction from $\text{MGNS}(d)$ to $\text{GNS}(2d)$, which will give us a $4d$ -approximation for $\text{MGNS}(d)$ using the algorithm in [19]. However, the reduction loses a lot of information about the topology of the underlying network.) For the case $d = \infty$, we derive an $O(\log n)$ -approximation for $\text{MGNS}(\infty)$ that matches the result of [19] for $\text{GNS}(\infty)$. Moreover, we give a polynomial time $\frac{3}{2}$ -approximation for $\text{MGNS}(1)$ on bipartite graphs.
2. We prove that for all $d \in \mathbb{N} \cup \{\infty\}$, it is \mathcal{APX} -hard to compute the optimum cost of $\text{GNS}(d)$ even if the graph is 3-regular and all costs and probability are uniform, thus ruling out the possibility of designing PTAS for the problem. We also show that, assuming the Unique Games Conjecture [13], we cannot obtain a $(\frac{3}{2} - \epsilon)$ -approximation for $\text{GNS}(d)$ in polynomial time. To our knowledge these are the first inapproximability results for $\text{GNS}(d)$. Since $\text{GNS}(d)$ is a special case of $\text{MGNS}(d)$, all the hardness results trivially apply to $\text{MGNS}(d)$.
3. We show that an optimal solution of $\text{MGNS}(d)$ can be found in polynomial time for every fixed $d \geq 1$ or $d = \infty$ if the underlying graph is a tree, and the infection cost and attack probability are both uniform. Our algorithm can handle all $d \leq O(\sqrt{\log n})$ in polynomial time on bounded-degree trees. Our algorithm also works for the case where budget constraints are put on the number of nodes and edges that can be secured and blocked respectively in a wanted solution. In particular, this settles an open question of [21] that asks whether there exists an efficient algorithm for the *minimum average contamination problem* on trees (which will be mentioned later in more detail). We remark that the tree structure, despite being special, has applications in hierarchically-organized networks such as company relationships.

Paper Organization. In the rest part of this section, we rigorously define our model and compare it with some previous work. In Section 2 we present

approximation algorithms for MGNS(d). Hardness of approximation results for MGNS(d) are given in Section 3. Section 4 copes with tree instances of MGNS(d). Finally, in Section 5 we conclude the whole paper and propose some open problems and future research directions.

1.1 Our Model for Infection Control

In this subsection we explain the mixed generalized network security model MGNS(d) in more detail, where $d \in \mathbb{N}^+ \cup \{\infty\}$ is a parameter that, intuitively, reflects the “degree of infectivity” within the network. Although we will describe our model in terms of preventing virus-spreading in computer networks, one should keep in mind that the model is capable of many other situations where we wish to minimize the propagation of undesirable things. Specifically, our model MGNS(d) comprises the following ingredients:

Contact Graph, Costs and Strategy. The contact graph is an undirected graph $G = (V, E)$, where $V = \{1, 2, \dots, n\}$ denotes the set of computers in a connected network, and $E \subseteq V^2$ specifies the underlying topology of the network. Thus, an edge $\{u, v\} \in E$ indicates that nodes (computers) u and v are directly connected, so that u can potentially affect v if it is infected by a computer virus or worm, and vice versa. For each $v \in V$, let C_v denote the *security cost* of v (for installing an anti-virus software on v), and L_v the *infection cost* of v (for recovering it from a virus attack). For each $e \in E$, let C'_e denote the *link-blocking cost* of e (for the lost caused by the removal of e). All the costs are non-negative. In a *strategy (solution)*, we need to decide on which nodes to install anti-virus softwares and which edges to block. A node with anti-virus software installed on it is called *secure*, and otherwise is called *insecure*. Similarly we have *blocked* and *unblocked* edges. A solution S is also identified with $V_S \cup E_S$, where $V_S \subseteq V$ is the set of secure nodes in S and $E_S \subseteq E$ is the set of blocked edges in S . The *attack graph* of a solution is the graph obtained from G by removing all secure nodes and blocked edges.

Infection Model and Social Cost. We assume that the virus is initiated at *exactly one* node chosen from V according to the *attack probability distribution* $\{w_v \mid v \in V\}$, where $\sum_{v \in V} w_v = 1$. Write $w(S) := \sum_{v \in S} w_v$ for $S \subseteq V$. A secure node will neither suffer from the virus nor transmit the virus to other nodes (although it can be chosen as the attacked node), whereas an insecure node, if chosen as the attacked node, will infect exactly those nodes at distance at most d from it in the attack graph (including itself). For a strategy S , let $V_S^{\leq d}(v)$ denote the set of nodes at distance at most d from v in the attack graph of S . Then the *social cost* of S (denoted by $cost(S)$) is defined as:

$$cost(S) = \underbrace{\sum_{v \in V_S} C_v}_{\text{cost for installing softwares}} + \underbrace{\sum_{e \in E_S} C'_e}_{\text{cost for blocking links}} + \sum_{v \in V \setminus V_S} \underbrace{L_v \cdot w(V_S^{\leq d}(v))}_{\text{expected cost for recovering } v \text{ from infection}} .$$

Goal. In the *centralized* setting of MGNS(d), the goal is to find a strategy with minimum social cost, or *social optimum*. We can also define the *decentralized*

(game-theoretic) model, in which the user needs to decide whether to install the anti-virus software on his/her computer and whether to disconnect some of the links with other users in the network. In this paper we concentrate on the centralized setting of $\text{MGNS}(d)$, while leaving explorations of the decentralized model to future work.

1.2 Related Work

As stated before, our model $\text{MGNS}(d)$ incorporates and generalizes several infection prevention models that have been studied recently. We list some problems considered in the literature that are either special cases of or related to the problem of computing the social optimum of $\text{MGNS}(d)$.

- Consider the instances of $\text{MGNS}(d)$ where $d = \infty$, $C'_e = \infty$ for all $e \in E$, all nodes have the same security cost C and infection cost L , and the attack probability distribution is uniform over nodes. When restricted on such instances, $\text{MGNS}(d)$ coincides with the model proposed by Aspnes, Chang and Yampolskiy [3], who gave an $O(\log^{1.5} n)$ -approximation for computing the social optimum, based on the sparsest cut algorithm of Arora, Rao and Vazirani [1]. The approximation ratio is subsequently improved to $O(\log n)$ independently by [5] and [19], which is also the currently best known result for this problem.
- Restricted on the instances where $C'_e = \infty$ for all $e \in E$ (i.e., all the edges should remain unblocked in any reasonable solution), our model is equivalent to the *generalized network security* model $\text{GNS}(d)$ introduced by Kumar et al. [19]. They present a $2d$ -approximation for computing the social optimum of $\text{GNS}(d)$ for all $d < \infty$ by rounding a natural linear program for the problem. This result is subsumed by our $(d + 1)$ -approximation for $\text{MGNS}(d)$. They also give an $O(\log n)$ -approximation for $\text{GNS}(\infty)$ based on a reduction to the *minimum weighted vertex multicut problem* [9], improving the $O(\log^{1.5} n)$ factor of [3] and matching the result independently obtained in [5].
- Under the case where $d = \infty$, $C_v = \infty$ for all $v \in V$, $w_v = 1/n$ for all $v \in V$, and both the infection costs and link-blocking costs are uniform, the problem of computing the social optimum of $\text{MGNS}(d)$ is similar to the *minimum average contamination problem* studied by Li and Tang [21], which originates from a (stochastic) link-blocking model initiated by Kimura, Saito and Motoda [16]. The difference between our setting and theirs is that they put a budget constraint K on the number of edges that can be removed from the network. In [21], a $(1 + \epsilon, O(\frac{\log n}{\epsilon}))$ -bicriteria approximation algorithm and a $(\frac{5}{3} - \epsilon)$ -inapproximability result are given for the minimum average contamination problem. Note that their problem is harder than ours (with an additional budget constraint) and thus their hardness factor is stronger than ours. However, they only consider the case $d = \infty$, while our hardness result applies to all d . Also, our polynomial-time algorithm for tree instances of $\text{MGNS}(d)$ holds for the budgeted case as well.

- Another related problem that has mainly been studied in the operations research forum is the *critical node problem* [2,4,7] defined as follows: given a node-weighted graph $G = (V, E)$, a connection cost $c(u, v)$ for each pair of nodes $\{u, v\} \in V^2$, and a parameter K , the goal is to find a subset of nodes whose total weight does not exceed K such that the total connection cost (counted for all connected pairs of nodes) is minimized. This problem is similar to MGNS(∞) with $C'_e = \infty$ for all $e \in E$ and $w_v = 1/n$ for all $v \in V$, but with additional budget constraints and more general cost functions. The problem is NP-complete on general graphs with unit costs and unit weights [2], and on trees with unit weights [7]. For the unit-cost case (which makes the problem fit in our model with $d = \infty$) in a tree of size n , Di Summa et al. [7] show that the problem is solvable in $O(n^7)$ time. Our polynomial-time algorithm for (budgeted) MGNS(d) on trees substantially generalizes their result to all fixed d .

2 Approximation Algorithm for MGNS(d)

In this section we concern with the computation of the social optimum of MGNS(d). As the problem is NP-hard, we focus on the perspective of approximation, and obtain the following results.

Theorem 1. *For any $d \geq 1$, there is a polynomial time $(d + 1)$ -approximation algorithm for computing the social optimum of MGNS(d). (Here d need not be a constant.)*

Theorem 2. *There is a polynomial time $O(\log n)$ -approximation for the social optimum of MGNS(∞).*

Theorem 3. *There is a polynomial time $\frac{3}{2}$ -approximation algorithm for computing the social optimum of MGNS(1) with bipartite contact graphs.*

We only prove Theorems 1 and 2 here. The proof of Theorem 3 will appear in the full version of this paper.

First consider the case $1 \leq d < \infty$. Let \mathcal{I} be an instance of MGNS(d) with contact graph $G = (V, E)$ where $V = \{1, 2, \dots, n\}$. If $C_i < w_i L_i$ for some $i \in V$, then clearly i should be secured in any optimum solution. Thus, we assume in what follows that $C_i \geq w_i L_i$ for all $i \in V$. We write an integer program to formulate the social optimum of \mathcal{I} . For each $k \in V \cup E$, let x_k be a binary variable that is 1 if and only if k is secure (or blocked, depending on whether k is a node or an edge). For a path p , let V_p and E_p denote the sets of nodes and edges on p , respectively. For all $1 \leq i < j \leq n$, let $P_{i,j}^d$ denote the collection of all simple paths from i to j of length at most d (note that $P_{i,j}^d$ can be empty and can also be of exponential size), and $y_{i,j}$ be a binary variable that is 1 if and only if there exists at least one path $p \in P_{i,j}^d$ on which all nodes are insecure and all edges are unblocked. Thus, $y_{i,j} = 1$ iff i and j can infect each other in the attack graph. Then the following integer program IP1 characterizes precisely the social optimum of \mathcal{I} :

$$\begin{aligned}
 \text{IP1: Min } & \sum_{i \in V} C_i x_i + \sum_{\{i,j\} \in E} C'_{\{i,j\}} x_{\{i,j\}} + \sum_{i \in V} L_i \left(w_i(1 - x_i) + \sum_{j \in V \setminus \{i\}} w_j y_{i,j} \right) \\
 \text{subject to: } & y_{i,j} + \sum_{k \in V_p \cup E_p} x_k \geq 1 \quad \forall 1 \leq i < j \leq n \text{ and } p \in P_{i,j}^d \\
 & y_{i,j} = y_{j,i} \quad \forall 1 \leq i < j \leq n \\
 & x_k \in \{0, 1\} \quad \forall k \in V \cup E \\
 & y_{i,j} \in \{0, 1\} \quad \forall 1 \leq i, j \leq n, i \neq j.
 \end{aligned}$$

We write $C'_i = C_i - w_i L_i$ for each $i \in V$ (with a little abuse of notation since C' is originally defined for edge costs), $L_{i,j} = w_i L_j + w_j L_i$ for all $1 \leq i < j \leq n$, and $C = \sum_{1 \leq i < j \leq n} w_i L_i$. Note that $C'_i \geq 0$ for all $i \in V$ by our assumption before. Rearranging terms, unifying the first two summations, and combining the occurrences of $y_{i,j}$ and $y_{j,i}$ in the objective function of IP1, we get a simpler yet equivalent formulation IP2 as follows:

$$\begin{aligned}
 \text{IP2: Min } & \sum_{k \in V \cup E} C'_k x_k + \sum_{1 \leq i < j \leq n} L_{i,j} y_{i,j} + C \quad \text{subject to:} \\
 & y_{i,j} + \sum_{k \in V_p \cup E_p} x_k \geq 1 \quad \forall 1 \leq i < j \leq n \text{ and } p \in P_{i,j}^d \\
 & x_k \in \{0, 1\} \quad \forall k \in V \cup E \\
 & y_{i,j} \in \{0, 1\} \quad \forall 1 \leq i < j \leq n.
 \end{aligned}$$

Observe that IP2, with the constant part C discarded, can be regarded as an instance of the weighted set cover problem when treating the length-at-most- d paths as the elements to be covered. When d is fixed, the instance of set cover is constructible in polynomial time. Also, in this set cover instance, every element appears in at most $2d + 2$ sets, because each constraint in IP2 involves at most $2d + 2$ variables (note that each $p \in P_{i,j}^d$ consists of at most $d + 1$ vertices and d edges). Therefore, a polynomial time $(2d + 2)$ -approximation exists for IP2 (see, e.g., [11]) and thus also for MGNS(d). Notice that, by reducing the problem to set cover, we can only handle constant d , and cannot hope for a poly-time $(2d + 2 - \epsilon)$ -approximation due to the $(k - \epsilon)$ -hardness of k -uniform hypergraph vertex cover [14], assuming the Unique Games Conjecture [13].

We next show that we can obtain an approximation factor of $d + 1$ for all d (not necessarily fixed) by utilizing the special structure of IP2, thus saving a factor of 2 from the set cover approach. To achieve this, we relax the last two constraints of IP2 to $x_k \geq 0$ and $y_{i,j} \geq 0$ respectively, and ignore the constant part C in the objective function. This gives us a linear programming relaxation (which might still have super-polynomial size) of the original instance, which we call LP. (We do not state LP explicitly since it is very similar to IP2.) Obviously, $OPT(LP) + C \leq OPT(IP2) = OPT(IP1)$, where $OPT(P)$ is the optimum objective value of the mathematical program P .

We now write the dual formulation of LP. Let $P^d = \cup_{1 \leq i < j \leq n} P_{i,j}^d$. For each $p \in P^d$, introduce a dual variable z_p , which corresponds to the constraint

$y_{i,j} + \sum_{k \in V_p \cup E_p} x_k \geq 1$ in LP (where i and j are the endpoints of p). The dual program DU can be written as follows:

$$\begin{aligned} \text{DU: Max } & \sum_{p \in P^d} z_p \quad \text{subject to:} \\ & \sum_{p \in P_{i,j}^d} z_p \leq L_{i,j} \quad \forall 1 \leq i < j \leq n \\ & \sum_{\substack{p \in P^d \\ k \in V_p \cup E_p}} z_p \leq C'_k \quad \forall k \in V \cup E \\ & z_p \geq 0 \quad \forall p \in P^d. \end{aligned}$$

By the strong duality theorem, $OPT(DU) = OPT(LP)$. We now find a solution to IP2 by Algorithm 1, which basically consists of a primal-dual procedure and a “pruning” phase. Since the number of variables in DU can be super-polynomial in n for non-constant d , the naïve implementation of Algorithm 1 may not run in polynomial time. Nevertheless, we will show later that the running time can be reduced to $n^{O(1)}$ regardless of d ; stating the algorithm in its current form is just to simplify the analysis of its performance guarantee. Let S denote the solution to IP2 returned by Algorithm 1, and $Z = \{z_p \mid p \in P^d\}$ be the solution to DU obtained in Algorithm 1 (which is not explicitly returned). Let $value(S)$ denote the objective value of the solution S .

Algorithm 1. Constructing a feasible solution for IP2

- 1: $x_k \leftarrow 0, \forall k \in V \cup E; y_{i,j} \leftarrow 0, \forall 1 \leq i < j \leq n$.
 - 2: $z_p \leftarrow 0, \forall p \in P^d$; also, set all z_p to be “unfrozen.”
 - 3: **while** there are still unfrozen variables **do**
 - 4: Choose any unfrozen variable, say z_p , that appears in some constraint of DU. Raise the value of z_p until some constraint in DU, say c , becomes tight. (Pick an arbitrary one if there are more than one tight constraints.)
 - 5: **if** c is “ $\sum_{p \in P^d: k \in V_p \cup E_p} z_p \leq C'_k$ ” for some $k \in V \cup E$ **then**
 - 6: $x_k \leftarrow 1$
 - 7: **else if** c is “ $\sum_{p \in P_{i,j}^d} z_p \leq L_{i,j}$ ” for some $1 \leq i < j \leq n$ **then**
 - 8: $y_{i,j} \leftarrow 1$
 - 9: **end if**
 - 10: Freeze all variables that occur in some (newly appeared) tight constraint.
 - 11: **end while**
 - 12: **for** all $1 \leq i < j \leq n$ **do**
 - 13: **if** $x_i = 1$ or $x_j = 1$ **then**
 - 14: $y_{i,j} \leftarrow 0; x_{\{i,j\}} \leftarrow 0$ if $\{i, j\} \in E$.
 - 15: **end if**
 - 16: **end for**
 - 17: **return** $\{x_k \mid k \in V \cup E\} \cup \{y_{i,j} \mid 1 \leq i < j \leq n\}$.
-

Lemma 1. *Z is a feasible solution to DU, and S is a feasible solution to IP2.*

The proof of Lemma 1 is easy and thus omitted.

Lemma 2. *value(S) ≤ (d + 1)OPT(IP2) .*

Proof. For each variable v of IP2, let $c(v)$ denote the constraint in DU that corresponds to v . Call a constraint $c(v)$ *active* if $v = 1$ in the solution S . By Line 4 of Algorithm 1, every active constraint $c(v)$ (say) is tight, and hence the contribution of this v to $value(S)$ (which is the coefficient of v in the objective function of IP2) equals to the sum of z_p 's contained in $c(v)$. Therefore, $value(S) = \sum_{p \in P^d} t_p z_p$, where t_p is the number of active constraints containing z_p .

Now fix an arbitrary $p = (i_0, i_1, \dots, i_t) \in P^d, t \leq d$. The set of constraints in which z_p appears is $\{c(y_{i_0, i_t})\} \cup \{c(x_{i_j}) \mid 0 \leq j \leq t\} \cup \{c(x_{\{i_j, i_{j+1}\}}) \mid 0 \leq j \leq t-1\}$, which can be partitioned into the following $t + 1$ subsets:

$$\{c(x_{i_0}), c(x_{\{i_0, i_1\}})\}, \{c(x_{i_1}), c(x_{\{i_1, i_2\}})\}, \dots, \{c(x_{i_{t-1}}), c(x_{\{i_{t-1}, i_t\}})\}, \{c(x_{i_t}), c(y_{i_0, i_t})\}.$$

Due to the function of the FOR loop, at most one constraint from each subset is active. Thus z_p appears in at most $t + 1 \leq d + 1$ active constraints. Recalling that the objective function of IP2 embraces an additional part C , we have

$$\begin{aligned} value(S) &\leq C + (d + 1) \sum_{p \in P^d} z_p \leq C + (d + 1)OPT(DU) \\ &= C + (d + 1)OPT(LP) \leq C + (d + 1)(OPT(IP2) - C) \\ &\leq (d + 1)OPT(IP2), \end{aligned}$$

completing the proof of Lemma 2. □

Lemmas 1 and 2 ensure that S is a $(d + 1)$ -approximate solution to IP2. We next explain how to make Algorithm 1 run in poly-time for all d . Consider the following two operations:

- (1) Find an unfrozen variable of DU if there exists at least one.
- (2) Given a variable z_p , find all the constraints in DU that contain z_p .

Lemma 3. *If operations (1) and (2) can be done in polynomial time, then Algorithm 1 can be implemented to run in polynomial time.*

Proof. Suppose (1) and (2) can be done in polynomial time. Since DU has at most $\binom{n}{2} + n \leq n^2$ constraints and each time only one variable raises its value, we can keep the current LHS and RHS values of each constraint, and are thus able to know which constraints are tight. Hence Line 10 can be realized implicitly since a variable is frozen iff it appears in some tight constraint. To implement Line 4, we first apply (1) to find an unfrozen variable (say z_p) if there exists one, and then use (2) to find a constraint containing z_p that has the smallest difference between RHS and LHS values; this difference is exactly the amount that z_p can be raised. The other steps in Algorithm 1 can clearly be implemented to run in poly-time. The lemma is thus proved. □

Lemma 4. *We can accomplish (1) and (2) in polynomial time.*

Proof. We use $c(v)$ to denote the constraint in DU that corresponds to the variable v of IP2. First note that (2) is easy to implement: For each variable z_p where p has endpoints i and j , z_p appears exactly in the constraints corresponding to $y_{i,j}$ or x_k for some $k \in V_p \cup E_p$. Thus we focus on (1). As shown in the proof of Lemma 3, we know the set of tight constraints in DU, and a variable is unfrozen if and only if it does not appear in any tight constraint. For $p \in P^d$, the variable z_p does not appear in $c(x_k)$ (where $k \in V \cup E$) iff $k \notin V_p \cup E_p$, and z_p does not appear in $c(y_{i,j})$ (where $1 \leq i < j \leq n$) iff p is not a path between i and j . We do the following: Construct a graph G' from G by deleting all $k \in V \cup E$ from G for which $c(x_k)$ is tight. Then, for every $1 \leq i < j \leq n$ such that $c(y_{i,j})$ is not tight, check whether there exists a path p from i to j in G' of length at most d ; if so, then the corresponding variable z_p must be unfrozen due to our previous analysis. Also, by this procedure we will find an unfrozen variable if there exists at least one. Clearly this process can be finished in polynomial time. \square

Now Theorem 1 follows directly from Lemmas 1, 2, 3 and 4.

We next turn to the case $d = \infty$ and prove Theorem 2. We reduce MGNS(∞) to GNS(∞) as follows: Construct a graph G' by subdividing each edge $e \in E$ with a new vertex v_e . Let $w(v_e) = 0, C_{v_e} = C'_e$ and $L_{v_e} = 0$ for all $e \in E$. It is easy to argue that the problem of finding the social optimum of GNS(∞) on this new instance is equivalent to that of MGNS(∞) on the original one. Now, applying the poly-time approximation algorithm for GNS(∞) given in [19], we get a solution for MGNS(∞) with approximation ratio $O(\log |V(G')|) = O(\log n)$. This finishes the proof of Theorem 2.

We remark that a similar reduction can reduce an instance of MGNS(d) to that of GNS($2d$). Using the approximation algorithm in [19], we obtain a solution for MGNS(d) with approximation factor $4d$, which is nearly four times larger than the ratio guaranteed by Theorem 1. This is in part due to the fact that such a reduction loses some information of the graph topology, which is important to our algorithm.

3 Hardness of Approximation for GNS(d)

In this section we present inapproximability results for GNS(d), a special case of our model MGNS(d). Thus, all the hardness results trivially apply to MGNS(d). The proof of the following two theorems will appear in the full version of this paper.

Theorem 4. *For every $d \in \mathbb{N} \cup \{\infty\}$, computing the social optimum of GNS(d) is APX-hard, even if the contact graph is 3-regular and all types of costs as well as the attack probability distribution are uniform.*

Theorem 5. *Assuming Unique Games Conjecture, for any $d \in \mathbb{N} \cup \{\infty\}$ and any fixed $\epsilon > 0$, we cannot approximate the social optimum of GNS(d) to a factor of $\frac{3}{2} - \epsilon$ in polynomial time.*

4 Polynomial Algorithm for MGNS(d) on Trees

In this section we consider a special class of instances of MGNS(d), in which the underlying contact graph of the instance is a tree, and the infection cost and attack probability are both uniform. Our main results are as follows, whose rigorous proofs will appear in the full version of this paper.

Theorem 6. *For every fixed $d \geq 1$ or $d = \infty$, we can find in polynomial time an optimal solution of a tree-instance of MGNS(d) with uniform infection cost and attack probability, even if there are budget constraints, i.e., given two integers K and K' , a solution can secure at most K nodes and block at most K' edges.*

Theorem 7. *For all $d \leq O(\sqrt{\log n})$, we can find in polynomial time an optimal solution to (budgeted) MGNS(d) if the instance has uniform infection cost and attack probability, and its contact graph is a tree of bounded degree.*

Theorem 6 in particular settles an open problem from [21] that asks if there is a polynomial time algorithm for the minimum average contamination problem, which corresponds to the special case of budgeted MGNS(d) on trees where every node has security cost ∞ and all other costs as well as the attack probability distribution are uniform.

5 Conclusions and Future Research

We propose in this paper the mixed generalized network security model MGNS(d), which generalizes several other models for infection control. We present approximation and inapproximability results for the problem of computing the optimum solution of MGNS(d), and exact polynomial-time algorithms for tree instances with uniform infection cost and attack probability distribution. Some of our results lead immediately to improvements upon the previously best known results achieved for some special cases of our model.

There are many interesting questions left that deserve further explorations. Regarding the optimization of social cost, a big open question is whether we can break the $O(\log n)$ factor for MGNS(∞) or GNS(∞), or there is a matching hardness of approximation result. Also for MGNS(d) where $d < \infty$, there remains a large gap between the upper bound of $d+1$ and the lower bound of $\frac{3}{2} - \epsilon$ on the approximation ratio. Another research issue is the formulation and investigation of the decentralized or game-theoretic counterpart of our model, where a user can decide whether to install an anti-virus software, and might also be able to block some of the links to other users. Finally, incorporating other propagation models (e.g., the independent cascade model, or the linear threshold model) into MGNS(d) may lead to more accurate modeling of some applications.

References

1. Arora, S., Rao, S., Vazirani, U.: Expander flows, geometric embeddings and graph partitioning. In: Proceedings of the 35th ACM STOC (2004)

2. Arulselvan, A., Commander, C.W., Eleftheriadou, L., Pardalos, P.M.: Detecting critical nodes in sparse graphs. *Comput. Oper. Res.* 36(7), 2193–2200 (2009)
3. Aspnes, J., Chang, K.L., Yampolskiy, A.: Inoculation strategies for victims of viruses and the sum-of-squares partition problem. *J. Comput. Syst. Sci.* 72(6), 1077–1093 (2005); Preliminary version in SODA 2005
4. Borgatti, S.: Identifying sets of key players in a social network. *Comput. Math. Org. Theory* 12, 21–34 (2006)
5. Chen, P.-A., David, M., Kempe, D.: Better vaccination strategies for better people. In: Proceedings of the 11th ACM EC (2010)
6. Chen, W., Wang, Y., Yang, S.: Efficient influence maximization in social networks. In: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2009)
7. Di Summa, M., Grosso, A., Locatelli, M.: Complexity of the critical node problem over trees. *Comput. Oper. Res.* 38(12), 1766–1774 (2011)
8. Domingos, P., Richardson, M.: Mining the network value of customers. In: Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2001)
9. Garg, N., Vazirani, V.V., Yannakakis, M.: Approximate max-flow min-(multi)cut theorems and their applications. *SIAM J. Comput.* 25(2), 235–251 (1993); Preliminary version in STOC 1993
10. Goyal, A., Bonchi, F., Lakshmanan, L.V.S.: Learning influence probabilities in social networks. In: Proceedings of the 3rd ACM WSDM (2010)
11. Halperin, E.: Improved approximation algorithms for the vertex cover problem in graphs and hypergraph. *SIAM J. Comput.* 31(5), 1608–1623 (2002)
12. Kempe, D., Kleinberg, J., Tardos, É.: Maximizing the spread of influence through a social network. In: Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2003)
13. Khot, S.: On the power of unique 2-prover 1-round games. In: Proceedings of the 34th ACM STOC (2002)
14. Khot, S., Regev, O.: Vertex cover might be hard to approximate to within $2 - \epsilon$. *J. Comput. Syst. Sci.* 74(3), 335–349 (2003); Preliminary version in CCC 2003
15. Kimura, M., Saito, K.: Tractable Models for Information Diffusion in Social Networks. In: Fürnkranz, J., Scheffer, T., Spiliopoulou, M. (eds.) PKDD 2006. LNCS (LNAI), vol. 4213, pp. 259–271. Springer, Heidelberg (2006)
16. Kimura, M., Saito, K., Motoda, H.: Minimizing the spread of contamination by blocking links in a network. In: Proceedings of the 23rd AAAI Conference on Artificial Intelligence (2008)
17. Kimura, M., Saito, K., Motoda, H.: Solving the Contamination Minimization Problem on Networks for the Linear Threshold Model. In: Ho, T.-B., Zhou, Z.-H. (eds.) PRICAI 2008. LNCS (LNAI), vol. 5351, pp. 977–984. Springer, Heidelberg (2008)
18. Kimura, M., Saito, K., Motoda, H.: Blocking links to minimize contamination spread in a social network. *ACM Trans. Knowl. Discov. Data.* 3(2) (2009)
19. Anil Kumar, V.S., Rajaraman, R., Sun, Z., Sundaram, R.: Existence theorems and approximation algorithms for generalized network security games. In: Proceedings of the 30th ICDCS (2010)
20. Leskovec, J., Krause, A., Guestrin, C., Faloutsos, C., VanBriesen, J., Glance, N.S.: Cost-effective outbreak detection in networks. In: Proceedings of the 13th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (2007)
21. Li, A., Tang, L.: The Complexity and Approximability of Minimum Contamination Problems. In: Ogihara, M., Tarui, J. (eds.) TAMC 2011. LNCS, vol. 6648, pp. 298–307. Springer, Heidelberg (2011)