

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/67537> holds various files of this Leiden University dissertation.

Author: Maulana, A.

Title: Many objective optimization and complex network analysis

Issue Date: 2018-12-05

7

Immunization of Networks Using Genetic Algorithms and Multi-Objective Metaheuristics

7.1 • Introduction

The study of networks has received increased attention in recent years. The effective control and combating of epidemics, such as Ebola [47] or the Zika virus [37], is one major problem, where the discovery of algorithms for analyzing and controlling networks can make an impact.

This chapter will focus on immunization strategies that achieve a high *eigenvalue drop*. The eigenvalue drop is the drop of the maximum eigenvalue after removal of a subset of nodes from a network, represented as an adjacency matrix. The eigenvalue drop is an effective measure for the impact of an immunization strategy because the maximum eigenvalue is inversely proportional to the epidemic threshold which determines how fast a virus spreads in the network and how long it lingers in the network [14, 15].

The epidemiological model that is considered in this work is the susceptible-infected-susceptible model, in short, SIS model. Here a node in the network can be infected via a direct neighbor and after a time it can recover and is susceptible again. See Figure 7.2 for different epidemiological models. Immunization of nodes can be enforced by measures outside of the network, e.g., by controlling the node or by removing the node from the network. In this work we assume that an immunized node can no longer infect other nodes, nor can it get infected itself.

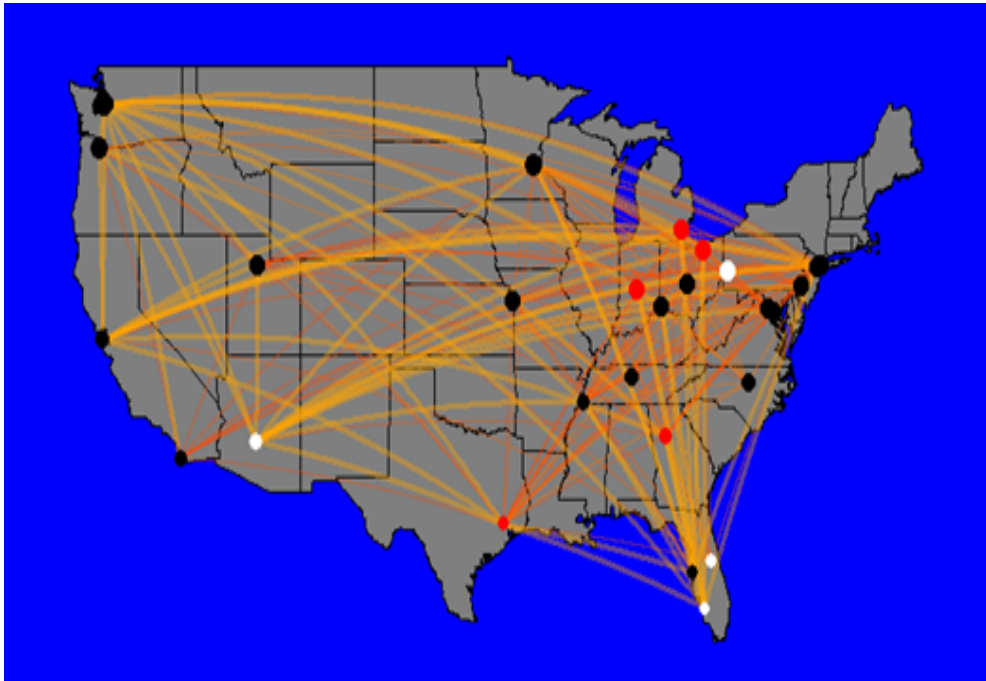


Figure 7.1 US flight network of major airports. The picture shows a snapshot of the spreading of a virus. The black nodes are susceptible, the white nodes are infected, and the red nodes are immunized. kateto.net/network-visualization

Consider for instance a network of airports connected by flights, such as the one provided in the US flights dataset shown in Figure 7.1. There might be some nodes already infected and we need to make it difficult for the virus to spread by controlling some major airports, e. g., by special bio-security checks or quarantining.

A network G will be represented as a pair (V, E) where V is a set of nodes $V = \{v_1, \dots, v_n\}$ and a set of edges $E \subseteq V \times V$. Vertices and edges can have weights and edge weights will be represented by a function $w_E : E \rightarrow \mathbb{R}^+$, and node weights by a function $w_V : V \rightarrow \mathbb{R}_0^+$. Given this, for a network we can alternatively use the adjacency matrix representation $A(V, E, w_E) \in \mathbb{R}^{n \times n}$ with $a_{ij} = 0$ if $(v_i, v_j) \notin E$ and $w((v_i, v_j))$ otherwise. The first or maximum eigenvalue of the graph will be denoted λ and the corresponding eigenvector with u . The components of this eigenvector, u_1, \dots, u_n , play a special role in this work and will be called the *eigen-scores* of the matrix.

Definition 7.1 Given a network G and a network G' , where G' is a subgraph of G with some nodes and their adjacent edges removed, the eigenvalue drop $\Delta\lambda$ is defined as

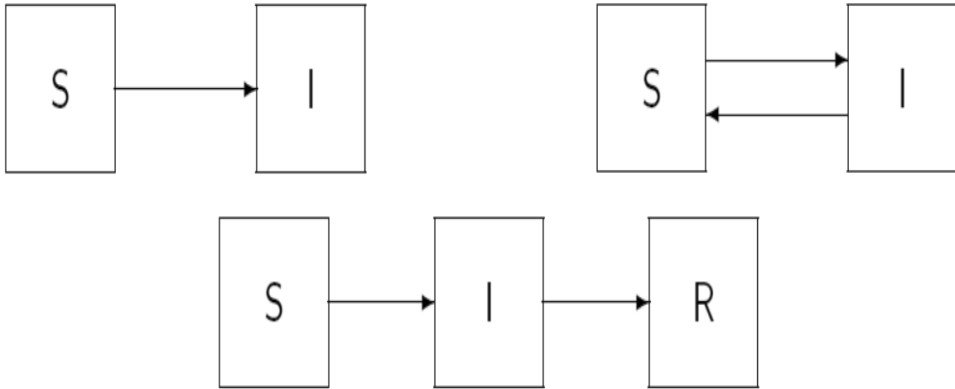


Figure 7.2 Three common models in epidemiology. In the SI model, nodes stay infected, once they got infected. In the SIS model, infected nodes can return into a susceptible state, and in the SIR model nodes are immunized after having recovered and can no longer infect neighboring nodes.

the difference between the maximum eigenvalue of the adjacency matrix of G and the maximum eigenvalue of the adjacency matrix of G' .

Definition 7.2 *The K -Node Immunization problem is the problem of finding a subset of k nodes to be removed from a network with n nodes, such that the eigenvalue drop is maximum.*

It has been shown in [15] that the decision problem that corresponds to the K -Node Immunization problem is NP-complete, and consequently the K -Node Immunization problem is NP-hard. Therefore, heuristic methods have been suggested in [15], most notably the Netshield Plus algorithm. This algorithm does not directly operate on the eigenvalue drop, but uses an approximation of it which is submodular and therefore lends itself to constructing an approximation algorithm. In brief, netshield seeks to maximize the following *Shield value* (S_v) function, which is closely correlated with the eigenvalue drop.

In this thesis, we propose an alternative approach to the k -node immunization problem based on genetic algorithms (Section 7.3.1) and compare results to Netshield Plus (Section 7.3.2). In the problem specific mutation operator, some of the ideas of Netshield will be adopted. Therefore, we will introduce this algorithm and the scoring function used by it briefly in Section 7.2. Moreover, a multi-objective generalization of

the k -node immunization problem is discussed. It introduces a cost term as a second objective (Section 7.4.1). First results on finding the Pareto front of this problem with multi-objective metaheuristics are presented in Section 7.4.2.

7.2 · Netshield Algorithm

Next, we will briefly introduce the Netshield algorithm. Some of the ideas of this algorithm will be useful in the design of the problem specific genetic algorithm. Moreover, the Netshield Plus algorithm, an improved version of the Netshield algorithm, will serve as a baseline algorithm in the benchmarking.

Let $G = (V, E)$ denote the original graph, and $G = (V', E')$ the graph after some nodes have been removed, and we define $S = V \setminus V'$. Moreover, A and A' denote the corresponding adjacency matrices. Then the Shield value (S_V) of S is defined as follows.

$$S_V(S) = \sum_{i \in S} 2\lambda(u_i)^2 - \sum_{i, j \in S} a_{ij}u_iu_j$$

Here, λ denotes the largest eigenvalue, u_i denotes the i -th component of the eigenvector that corresponds to the largest eigenvalue. It is also called the i -th eigen-score. The Shield value rewards dissimilarity between nodes, that is small a_{ij} , and a high eigen-score.

As opposed to the Netshield algorithm, the Netshield Plus algorithm [15] removes nodes in batches of b nodes each. After each batch, the largest eigenvalue and the corresponding eigen-scores are recomputed. This way the algorithm yields more accurate results, but due to multiple eigenvalue computations the computation time increases. Netshield Plus is therefore especially recommended for small or moderate size networks, as we discuss them in this chapter.

7.3 · Problem Specific Genetic Algorithm

7.3.1 · Discussion of the method

In this work we use a standard $(\mu + \mu)$ genetic algorithm (see, e.g., [60]) with scaled proportional selection (mating selection) and truncation selection (environmental selection). The genetic algorithm for the k subset selection problem uses problem specific mutation and crossover operators. The representation of solution candidates

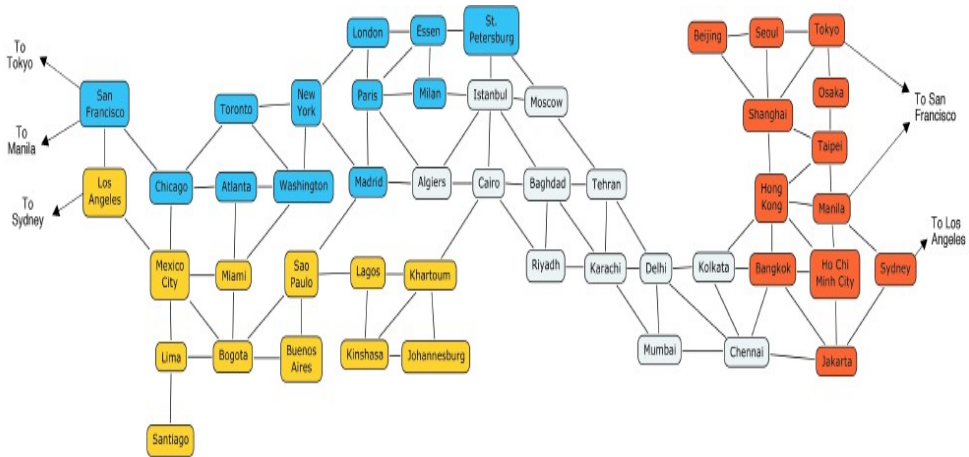


Figure 7.3 Network in the Pandemic game board
(from: <http://jhkimrpg.livejournal.com/78787.html>).

is not binary, as usual, but a problem specific representation for subset selection as it has been used in other contexts, too [61]. A solution is represented as k non duplicate integers in $[1, n] \subset \mathbb{N}$.

The mutation operator that was designed for this problem relies on two mechanisms:

- Firstly, in each mutation, an integer that is in the array is replaced by an integer in $[1, n]$ that is not in the array.
- Secondly, the algorithm works with two different mutation rates. For nodes with a top- k eigen-score, the probability of mutation is increased by a constant factor ≥ 0 , making it more likely to be selected for the set or discarded. This way it is hypothesized that the algorithm spends more time in exploring relevant parts of the graph. The multiplication factor will be denoted with ν .

Mutation is applied to each offspring individual. First, an integer in the array is selected proportionally to the mutation probabilities. Then an integer outside the array is selected proportional to the mutation probabilities. And then the node inside the array is replaced by the node outside the array. The genetic algorithm does not feature crossover, but we might consider the development of a problem specific crossover for future research.

7.3.2 · Comparison to Netshield Plus

For the empirical comparison of algorithms we will use five data sets on networks:

- Karate: A social network of friendships between 34 members of a karate club at a US university in the 1970s [62].
- Dolphins: It is a social network consisting of an undirected network of frequent associations between 62 dolphins in a community living off Doubtful Sound, New Zealand. [38]
- US Flights: This is a list of the most important Airports in the United States connected to other based on the exist are of connecting flights (edge) from one airport to the other airports.
- Pandemic: A cooperative board game with the goal to fight the outbreak of the virus. We used the graph that connects cities in the world as an example data set [35]. A picture of the Pandemic board is shown in Figure 7.3
- Conference Day 1: The social interactions of members of a conference on the first day. Taken from <http://www.sociopatterns.org/datasets/infectious-sociopatterns>.
- Conference Day 3: From the same data set as above, but for the third day.

The data sets US flights and Pandemic are most representative for the problem class. The other networks are added to gain more general insights into the algorithm behavior and reliability. Note that social interaction networks are also relevant in the spread of the virus, albeit control is less straightforward as compared to networks where nodes are assigned to places, such as US flights and the Pandemic board game network.

For the k -node immunization problem we used the Netshield Plus algorithm and parameters as described in [15]. For the genetic algorithm tests the following setting was applied: The number of function evaluations was 30000. Different mutation parameters were tested, with a value of $\nu \in \{1/n, 2/n, 3/n, 6/n, 1\}$, that is the mutation rate for the k components of u with the highest eigen-score. For all other nodes, the mutation probability was set to $1/n$, which is a recommended rate according to Bäck [4].

Regarding the single objective genetic algorithms, they were executed 20 times each, for $k = 3, 5$ and 10 on the Karate, Dolphins, US Flights, Pandemic, Conference Day 1 and Conference Day 3 networks. Table 7.1 shows results for single objective optimization of the eigenvalue drop. For assessing statistical significance we also provide box plots of our results in Figure 7.5, 7.6, 7.7 and Figure 7.6. We observe that GA_5, which represents the $(\mu + \mu)$ genetic algorithm that introduces a mass of

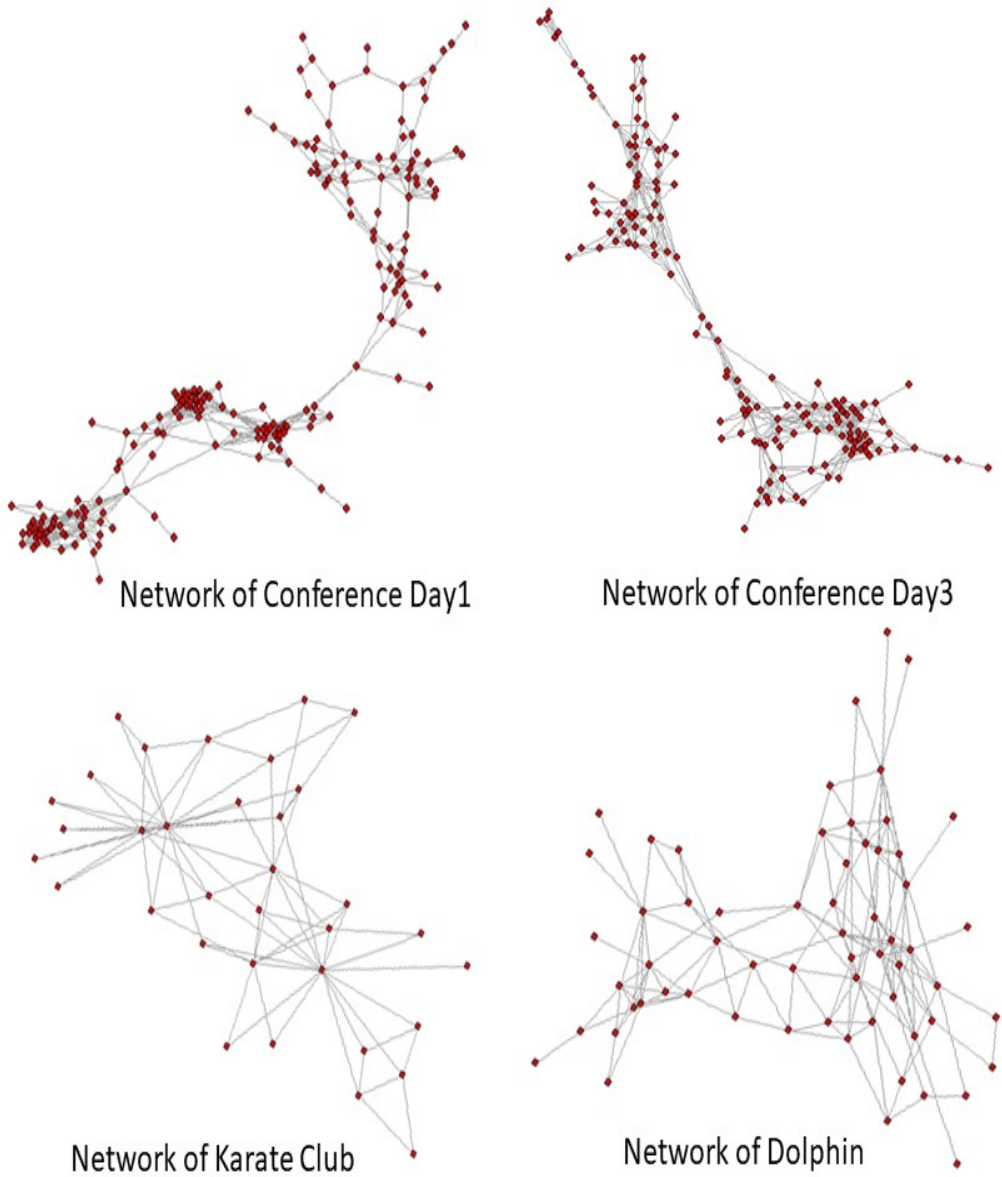


Figure 7.4 Depiction of network that we consider for our experiments consist of network of Conference Day1, Conference Day3, Karate Club and network of Dolphin

	<i>Network</i>	<i>GA_1</i>	<i>GA_2</i>	<i>GA_3</i>	<i>GA_4</i>	<i>GA_5</i>	<i>Netshield+</i>
<i>K = 5</i>	<i>karate</i>	4.1068	4.1068	4.1068	4.1068	4.1068	4.1068
	<i>Dolphins</i>	2.0812	2.0769	2.0807	2.0978	2.0978	2.0817
	<i>USA</i>	7.2043	7.2043	7.2043	7.2043	7.2043	7.2043
	<i>Pandemic</i>	0.9243	0.9419	0.9502	0.9502	0.9133	0.9556
	<i>Conf.day1</i>	3.0109	2.9583	3.0289	3.0455	3.0391	3.0638
	<i>Conf.day3</i>	17.670	17.671	17.669	17.669	17.610	3.8542
<i>K = 10</i>	<i>karate</i>	5.1077	5.1077	5.1077	5.1077	5.3115	5.3115
	<i>Dolphins</i>	2.9077	2.9230	2.9685	3.1575	3.2862	3.3997
	<i>USA</i>	11.690	11.809	11.922	12.177	12.608	12.608
	<i>Pandemic</i>	1.4201	1.4299	1.4490	1.5114	1.5215	1.4442
	<i>Conf.day1</i>	4.3853	4.3831	4.4207	4.6697	19.237	4.9121
	<i>Conf.day3</i>	17.659	17.664	17.658	17.659	17.658	5.6483

Table 7.1 Results of genetic algorithm and Netshield Plus comparisons.

5 to the k -highest eigen-score nodes, to be the best candidate. Although there is not a unanimously best genetic algorithm for the task, we consider our genetic algorithms to be a supplementary tool to Netshield/Netshield Plus, for medium-sized networks (≤ 200 nodes).

7.4 · Multi-Objective Node Immunization

In real-world scenarios, it is likely that multiple nodes need to be controlled or immunized, but it is typically not the case that the value of k is given a priori. Rather it is the case that the immunization of a node comes with a cost, which can differ from node to node. If a larger number of nodes is immunized the total cost would be approximately proportional to the cumulated cost of immunizing the single nodes. Let S denote the set of indexes of the immunized nodes and c_i denote the cost of immunization of node i , defined a priori. Then the *immunization cost* objective function can be defined as

$$C(S) = \sum_{i \in S} c_i \rightarrow \min$$

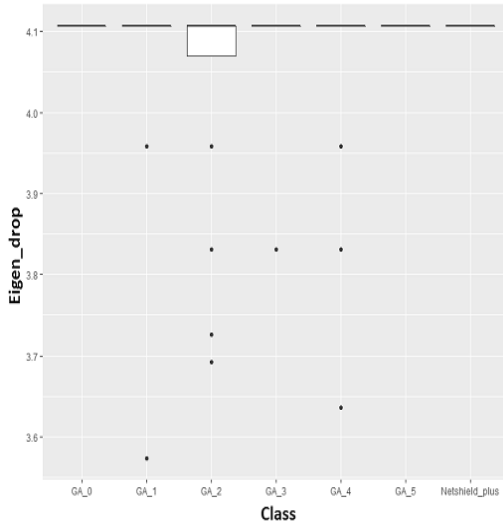
In multi-objective optimization, problems with two or more objectives are solved. In case of the node immunization problem the problem formulation reads as follows:

<i>ID</i>	<i>City</i>	<i>Population</i>	<i>ID</i>	<i>City</i>	<i>Population</i>
1	<i>SanFrancisco</i>	723724	25	<i>Beijing</i>	7602069
2	<i>Chicago</i>	2830144	26	<i>Seoul</i>	9860000
3	<i>Montreal</i>	3280123	27	<i>Tokyo</i>	8372440
4	<i>NewYork</i>	8124427	28	<i>Shanghai</i>	15017783
5	<i>Washington</i>	548359	29	<i>HongKong</i>	7347000
6	<i>Atlanta</i>	424096	30	<i>Taipei</i>	2491662
7	<i>Madrid</i>	3146804	31	<i>Osaka</i>	2590815
8	<i>London</i>	7489022	32	<i>Bangkok</i>	4935988
9	<i>Paris</i>	2141839	33	<i>HoChiMinhCity</i>	3496586
10	<i>Essen</i>	596204	34	<i>Manila</i>	10546511
11	<i>Milan</i>	1316218	35	<i>Jakarta</i>	8556798
12	<i>St.Petersburg</i>	4991000	36	<i>Sydney</i>	4444513
13	<i>Algiers</i>	2029936	37	<i>Khartoum</i>	2090001
14	<i>Istanbul</i>	10034830	38	<i>Johannesburg</i>	2091491
15	<i>Moscow</i>	10472629	39	<i>Kinshasa</i>	9464000
16	<i>Cairo</i>	7836243	40	<i>Lagos</i>	9020089
17	<i>Baghdad</i>	5753612	41	<i>SaoPaulo</i>	10059502
18	<i>Tehran</i>	7160094	42	<i>BuenosAires</i>	11595183
19	<i>Delhi</i>	11215130	43	<i>Santiago</i>	4893495
20	<i>Karachi</i>	11969284	44	<i>Lima</i>	7857121
21	<i>Riyadh</i>	4328067	45	<i>Bogota</i>	7235084
22	<i>Mumbai</i>	18410000	46	<i>MexicoCity</i>	8659409
23	<i>Chennai</i>	7088000	47	<i>LosAngeles</i>	3911500
24	<i>Kolkata</i>	4497000	48	<i>Miami</i>	386740

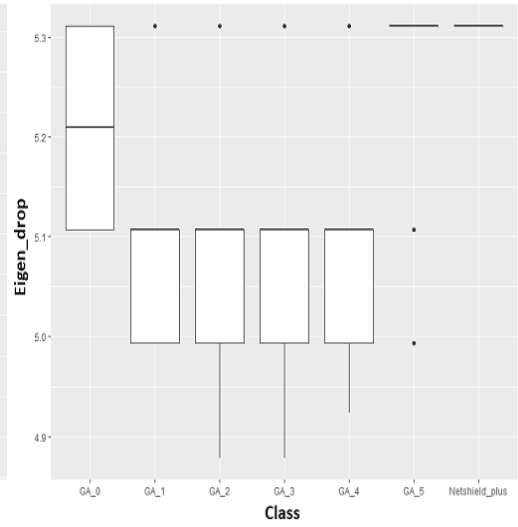
Table 7.2 Cost values for Pandemic network (proportional to city size).

<i>Label</i>	<i>Airport</i>	<i>Visits</i>
42	<i>Cincinnati/northernKentucky</i>	117
51	<i>DetroitMetropolitanWayneCounty</i>	126
71	<i>GeorgeBushIntercontinental</i>	90
81	<i>Hartsfield – jacksonAtlantaInternational</i>	102
85	<i>HopkinsInternational</i>	123
88	<i>IndianapolisInternational</i>	120
106	<i>KansasCityInternational Airport</i>	117
110	<i>LaGuardia</i>	123
131	<i>MemphisInternational</i>	105
137	<i>Minneapolis – St.PaulIntl</i>	135
153	<i>NashvilleInternational</i>	108
155	<i>NewarkLibertyInternational</i>	123
164	<i>OrlandoInternational</i>	84
169	<i>PhiladelphiaInternational</i>	120
172	<i>PittsburghInternational</i>	120
173	<i>PortColumbusIntl</i>	120
174	<i>PortlandInternational</i>	138
177	<i>Raleigh – durhamInternational Airport</i>	108
190	<i>RonaldReaganWashingtonNational Airport</i>	117
193	<i>SaltLakeCityInternational</i>	123
195	<i>SanDiegoInternational Airport</i>	99
196	<i>SanFranciscoInternational</i>	114
201	<i>Seattle – TacomaInternational</i>	141
204	<i>SkyHarborIntl</i>	99
206	<i>SouthwestFloridaReg</i>	81
214	<i>TampaInternational</i>	84
224	<i>WashingtonDullesInternational</i>	117

Table 7.3 Cost values for Pandemic network (proportional to city size).

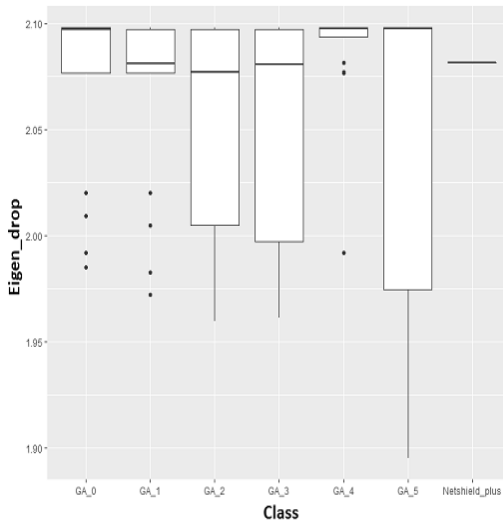


(Karate $k = 5$)

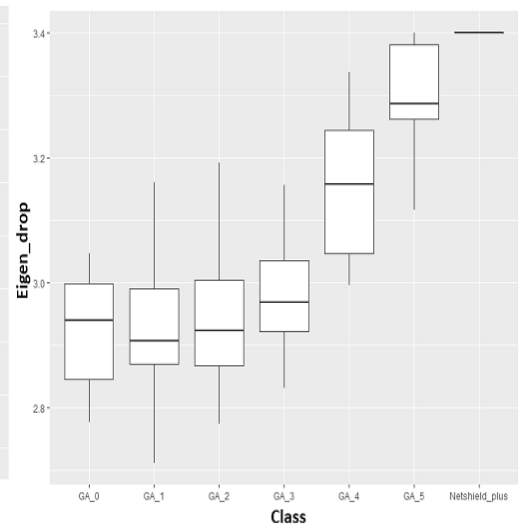


(Karate $k = 10$)

Figure 7.5 Results of genetic algorithm and Netshield Plus comparisons for the Karate network.



(Dolphins $k = 5$)



(Dolphins $k = 10$)

Figure 7.6 Results of genetic algorithm and Netshield Plus comparisons for the Dolphin network.

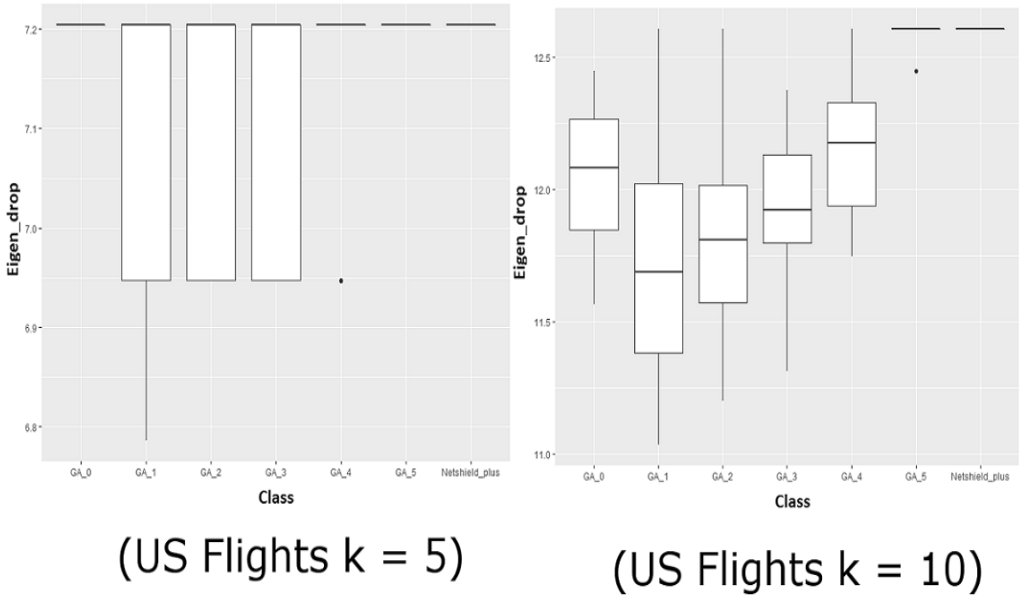


Figure 7.7 Results of genetic algorithm and Netshield Plus comparisons for the US flight network.

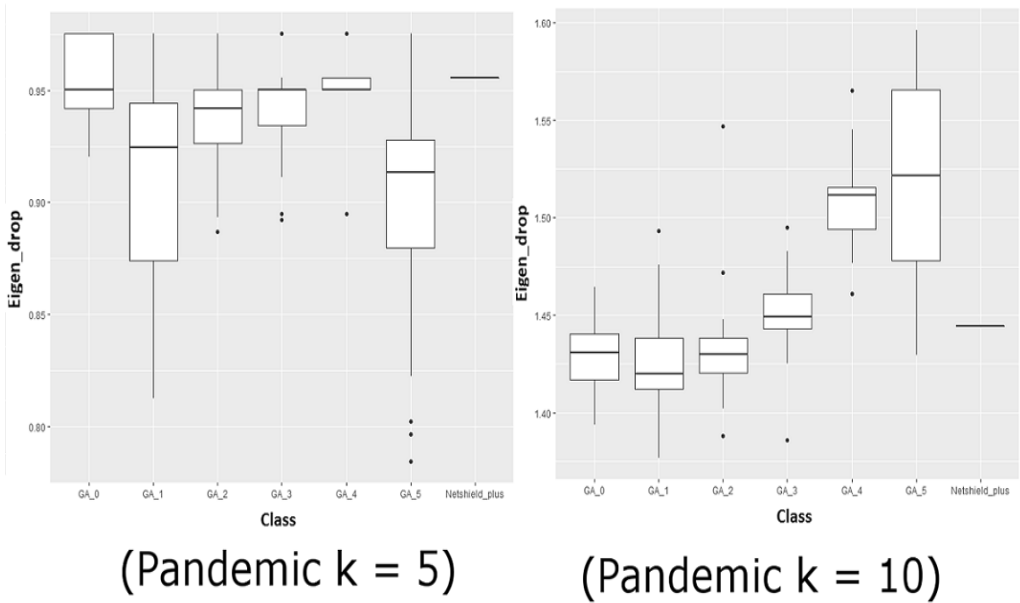


Figure 7.8 Results of genetic algorithm and Netshield Plus comparisons for the Pandemic network.

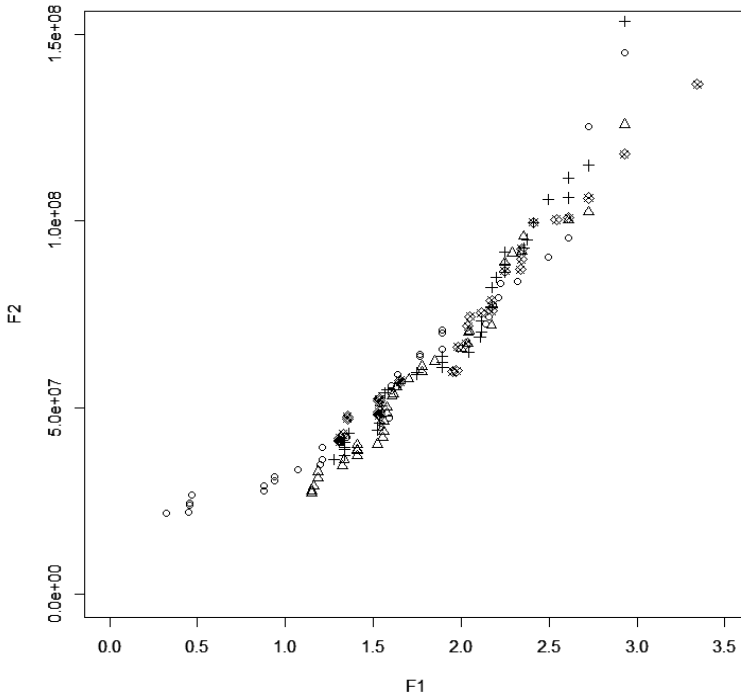


Figure 7.9 Pareto Front for the Pandemic Network found by the NSGA-II algorithm with 5 experiments.

$$f_1(S) = \lambda(S) \rightarrow \max \quad (7.1)$$

$$f_2(S) = C(S) \rightarrow \min \quad (7.2)$$

$$S \subseteq \{1, \dots, n\} \quad (7.3)$$

We are interested in the efficient set of this problem, that is the set: $\mathcal{S}_E = \{S \in \{1, \dots, n\} \mid \nexists S' \subseteq \{1, \dots, n\} : f_1(S') \geq f_1(S) \wedge f_2(S') < f_2(S) \vee f_1(S') > f_1(S) \wedge f_2(S') \leq f_2(S)\}$ and the Pareto front $\{(f_1(S), f_2(S))^T \mid S \in \mathcal{S}_E\}$.

7.4.1 · Multi-objective Metaheuristics

Two multi-objective evolutionary algorithms (MOEA, or EMOA) are considered as solvers: The first one is the non-dominated sorting genetic algorithm (NSGA-II) [17] and the second one is the S-metric selection algorithm (SMS-EMOA) [8].

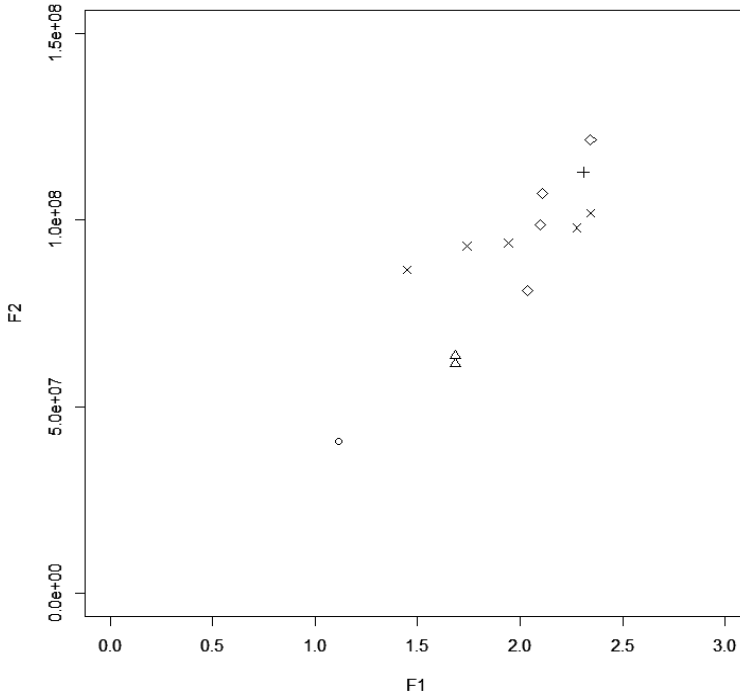


Figure 7.10 Pareto Front for the Pandemic Network found by the SMS-EMOA algorithm with five experiments.

The implementations of SMS-EMOA and NSGA-II in R featured by Bossek’s `ecr` package was used in this work. The representation of a subset is chosen to be a bit vector b in \mathbb{B}^n , where $b_i = 1$ means the node is selected to be removed/quarantined and $b_i = 0$ means the node is not selected, for $i = 1, \dots, n$. As recombination operator, one point crossover is used. For all bits, we used $p_m = 1/n$ as the mutation probability. The reason for this mutation rate is that, in contrast to the single objective genetic algorithms we discussed, here we do not know a priori the number of nodes to remove/quarantine. That is, we do not specify a subset cardinality. As a consequence, the algorithm should not try to explore a particular direction of the search space (bias introduced from the mutation operator), but rather present to the decision makers a complete picture of their possible choices. For example, quarantining 10 less-important (in terms of eigen-score) airports could be more beneficial in terms of cost, than quarantining 1 important (in terms of eigen-score) airport.

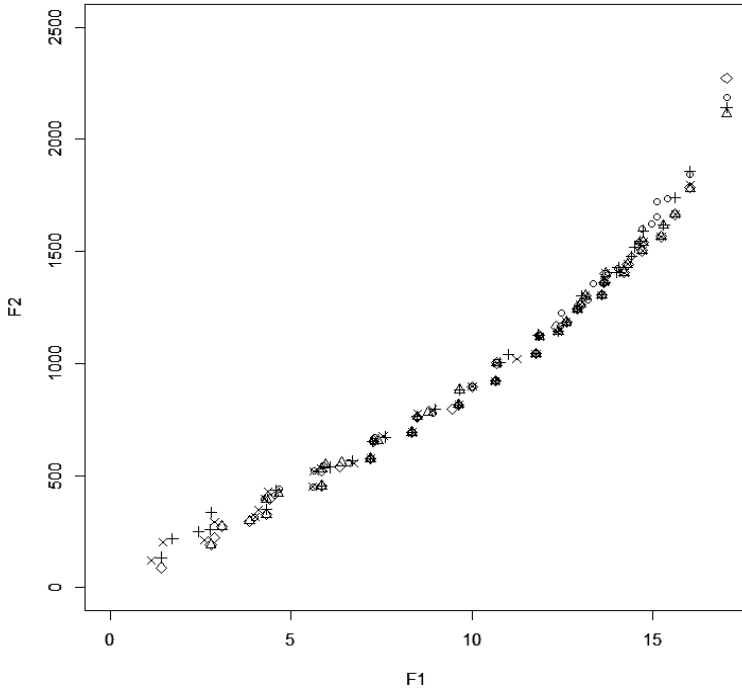


Figure 7.11 Pareto Front for the USA Flight Network found by the NSGA-II algorithm with 5 experiments.

7.4.2 · Empirical Results

The Pandemic and the US flights networks serve as examples for computing the Pareto fronts and efficient sets. In case of the Pandemic network, the size of the cities was used as a cost, assuming that it is more difficult to immunize larger cities. In case of the US flights network, the size of the airport (number of visits) was taken into account. The cost values are tabulated in Table 7.2 (Pandemic) and Table 7.3 (US flight). While we aimed for realistic problem settings, we would like to note that in order to plan effective real-world immunization more modeling is needed, including social interactions, geographic environment, and various other factors. Here, we merely focus on the network aspects of the problem. Each algorithm for the multi-objective optimization was run 5 times, producing 5 Pareto front approximations. Results for the Pandemic Network are shown in Figure 7.9 and Figure 7.10. Results for the Pandemic are shown in Figure 7.11 and Figure 7.12. The Pareto front looks near linear. Overall the

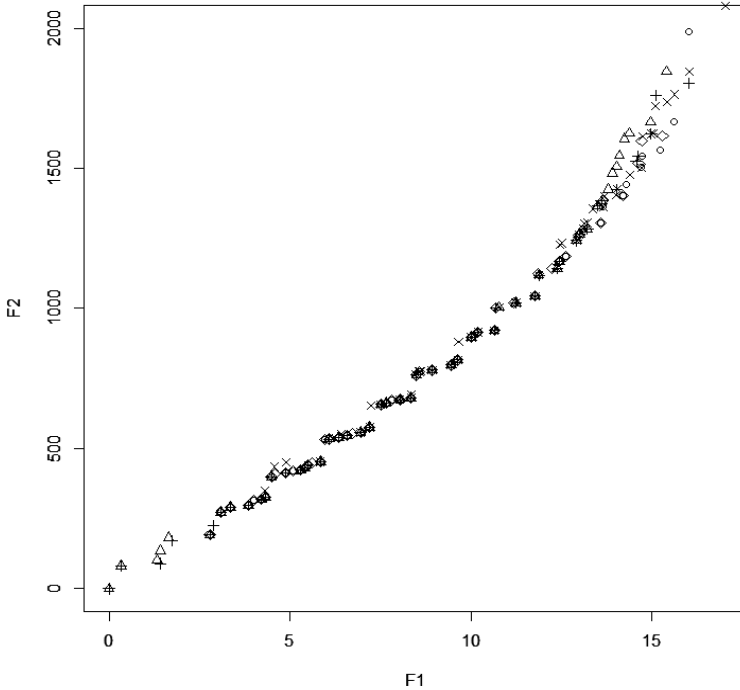


Figure 7.12 Pareto Front for USA Flight Network found by the SMS-EMOA algorithm with 5 experiments.

NSGA-II algorithm obtained better results and displayed a more robust performance than the indicator based SMS-EMOA on this problem. It is also interesting to note that the Pareto front looks near linear, which might be explained by the fact that big nodes (larger cities or, respectively, airports) are at the same time costly and important for immunization. For the US Flights network, a knee region can be identified.

7.5 · Summary

This chapter discusses network immunization techniques based on a heuristic method using genetic algorithms technique. Compared to Netshield Plus, the results show that the genetic algorithm often performs better, sometimes significantly better, in solving the k -node immunization problem. Netshield Plus is a fast heuristic and produces in many cases good results. Based on our findings, a strategy could be recommended that, if time is available, uses not only Netshield Plus but also a problem specific genetic algorithm to make it more likely that the best solution for the edge drop objective is

not overlooked.

In order to achieve good results, specific adaptations turned out to be very useful. An idea that works well is to use eigen-score values in order to adjust the mutation probabilities. This way the search is more focused on the part of the search space that is more likely to be relevant to solving the problem. We should also emphasize here that the supplementary use of a problem specific genetic algorithm has the advantage of calculating the actual eigen-drop, rather than an approximation of it. This can be useful for moderately sized networks. However, in large networks, the computational cost increases, since the algorithm eigen decomposes larger adjacency matrices.

First results were also obtained on a multi-objective formulation of the node immunization problem. We discuss the formulation where the total cost of immunization is one objective and the drop of the eigenvalue is the second objective. Two different meta-heuristics are applied to solve this problem and they widely agree with the results and show robust performance.¹

¹We remark that the source code of the algorithms and the network datasets are available by the authors on the research groups web page <http://moda.liacs.nl>.

