

Near-Optimal Blacklisting

Christos Dimitrakakis

EPFL

Lausanne, Switzerland

Email: christos.dimitrakakis@epfl.ch

Aikaterini Mitrokotsa

Chalmers University of Technology

Gothenburg, Sweden

Email: mitrokatkm@gmail.com

Abstract—Many applications involve agents sharing a resource, such as networks or services. When agents are honest, the system functions well and there is a net profit. Unfortunately, some agents may be malicious, but it may be hard to detect them. We consider the intrusion *response* problem of how to permanently *blacklist* agents, in order to maximise expected profit. This is not trivial, as blacklisting may erroneously expel honest agents. Conversely, while we gain information by allowing an agent to remain, we may incur a cost due to malicious behaviour. We present an efficient algorithm (HiPER) for making near-optimal decisions for this problem. Additionally, we derive three algorithms by reducing the problem to a Markov decision process (MDP). Theoretically, we show that HiPER is near-optimal. Experimentally, its performance is close to that of the full MDP solution, when the (stronger) requirements of the latter are met.

I. INTRODUCTION

We consider the *decision making* problem of blacklisting potentially malicious nodes or agents that share a resource or network based on partial information. As motivation, consider a communication network which is monitored by a network management system. Nodes can be of one of two types: malicious (e.g. dropping or corrupting packets, creating undue congestion), or honest. At each time-step (e.g. reporting period), we get a set of readings, giving some information about the behaviour of each node during that period. We want to optimally decide whether to blacklist a node, or maintain it in the system for one more time-step.

In order to make the problem non-trivial, we assume that for every honest node in the network, we have some fixed tangible gain at each time period. This would be the case, if all participation was done through a subscription model as in internet service providers (ISPs). On the other hand, we incur a (hidden) cost for each malicious node that participates. Thus, it is in our interests to kick out malicious nodes as soon as possible, but never to expel honest ones.

We should emphasize that this is *not* an intrusion detection problem. In fact, the readings that we obtain for each node could be seen as the output of some intrusion detection system (IDS). Rather, we are more concerned about the decision making aspect: what is the *optimal response* to the IDS outputs, given assumptions about the cost of malicious behaviour?

This setting of keeping suspicious nodes in the network until we become more certain about their type appears in many applications such as: 1) blacklisting clients of an ISP 2) shutting down malware-infected hosts in an internal network

3) expelling selfish nodes from a peer-to-peer network. In all of the above cases, any single piece of information is not enough to condemn a node to blacklisting. Rather, a sufficient amount of statistics has to be collected before we are sure that removing a node is more beneficial than keeping it. In this paper, we propose and consider a number of algorithms for tackling this problem in a general setting.

More precisely, our first contribution is a decision-theoretic approach based on distribution-free high probability bounds. The bounds require very little prior information and can be used to trade off the cost of removing honest nodes with that of keeping malicious nodes in the network for too long. We prove that this High Probability Efficient Response algorithm (HiPER) has low *worst-case expected loss* relative to an oracle which knows the type of every node.

Our second contribution is a set of Bayesian decision-theoretic approaches that we derive by formalising the problem as a Markov decision process (MDP). These require some further assumptions. In particular, it is necessary to fully specify a structure and prior parameters for the underlying statistical model. In addition, making optimal decisions according to such models is computationally intractable. Consequently, we consider some approximate algorithms. Of these, an optimistic approximation has similar performance to that of HiPER, while a finite lookahead approximation has increased performance, at the cost of additional computation.

The paper is organised as follows. In the remainder of this section we give some background, present related work and our contributions. Section II introduces notation while Section III specifies the loss model. Section IV presents the proposed HiPER algorithm as well as the bounds on the *worst-case expected loss*. Section V describes the decision-theoretic approaches which model the problem as an MDP and are used in the performance comparisons with the HiPER algorithm while Section VI describes the evaluation experiments. Finally, Section VII concludes the paper. The appendix provides proofs of technical lemmas and some useful auxiliary results.

A. Background

The problem we consider falls within decision theory. In particular, the scenario we investigate can be reduced to the *optimal stopping* problem [7], which can be modelled as an MDP [7] or as a (potentially unknown) partially observable MDP (POMDP) [18].

More precisely, in our setting, the nodes can be one of two types: honest or malicious. However, we initially start out without knowing what type each node is. Consequently, we must gather data (observations) to reduce our uncertainty about their types. Unfortunately, we can only do so while a node remains within the network. However, the longer we maintain a malicious node in the network, the more loss we incur. Conversely, once we remove an honest node, we will obtain no more profit from it. So, the problem can be reduced to deciding at what time, or under which conditions, to remove a given node from the network, if at all. Thus, our scenario can be seen as a type of *optimal stopping* problem.

The stopping problem has been extensively studied in general [7], while partial monitoring games in general have also received a lot of attention recently [6]. However, to the best of our knowledge, the general hidden reward stopping problem has not been previously studied in the literature. On the other hand, the specific application we consider can be seen as a type of *optimal intrusion response*.

Most of the previous research on intrusion response has concentrated on the POMDP formalism. Indicative publications are those by Zonouz *et al.* [22], Zan *et al.*, [20] and Zhang *et al.* [21], which have all proposed an intrusion response through modelling the process as a POMDP [18]. More precisely, Zonouz *et al.* [22] proposed a Response and Recovery Engine (RRE) based on a game-theoretic response strategy against adversaries modelled as opponents in a non-zero-sum, two-player Stackelberg stochastic game. In each step of the game RRE chooses the response actions using an approximate POMDP solver. More precisely, using the most likely state (MLS) [5] approximation, the POMDP is converted to a competitive Markov Decision Process (MDP), which is then solved using a look ahead search (i.e. approximate planning). Zhang *et al.* use the POMDP to integrate low level IDS alerts with high level system states, while Zan *et al.* [20] propose to solve the intrusion response problem as a factored POMDP model. Additionally, they decompose the POMDP into small sub-POMDPs and compute the response policy using the MLS approximation technique. However, in our case MLS as an approximation is too crude to be used, since it would essentially result in a completely random policy, as there are only two possible hidden states each node can be in. An entirely different approach, policy-gradient methods, is employed by [8] in the context of combating denial-of-service attacks in P2P networks. However, this approach requires observing the rewards, which are in fact hidden in our case.

B. Our contributions

Our first proposed algorithm relies on bounds which do not require knowledge of prior probabilities regarding the type of a node (honest or malicious) neither known distributions for the observations corresponding to honest or malicious nodes. We only need to know the mean of each of these distributions. Consequently, it is substantially more lightweight than MDP solvers, since we take decisions without performing explicit planning. Thus, it is more suitable for resource constrained

environments. We analyse the *expected loss* of this algorithm, and show that it is not significantly worse to that of an oracle which already knows each node's type.

Our second contribution is to derive three approximate MDP solvers. In contrast to previous work, in our scenario the reward is *never observed* by the algorithm.¹ This corresponds to reality, since we frequently do not know which nodes give us negative rewards.² Furthermore, two of our MDP algorithms are different from those previously employed in the intrusion response literature, as we forego the most-likely-state approximation commonly used in POMDP approaches. We first consider a myopic approximation.³ The second approach is a lightweight *optimistic* approximation that performs no planning, which is derived from upper bounds [9] on Bayesian decision making in unknown MDPs [11]. To our knowledge, this approach has not been used in similar problems before. Finally, we consider online planning with finite lookahead [15], [7]. This approach takes decisions which consider the impact of all our possible future actions up to some horizon. This approach has been employed in other applications such as dialogue modelling [4], autonomous underwater vehicle mapping [16], preference elicitation [2] and sensor scheduling [12] in wireless sensor networks.

II. PRELIMINARIES

We consider a set of nodes, which can be either honest or malicious. We assume there is a reliable way to obtain statistics from each node, such as an IDS that gives us a numerical score for each node. We denote by \mathcal{Q} the set of all malicious nodes and by \mathcal{U} the set of all honest nodes. We consider that there is an entity \mathcal{E} (for instance an Internet Service Provider (ISP) or a network administrator) who gains some reward $g_{\mathcal{U}}$ for each moment that an honest node remains in the network and has a cost $\ell_{\mathcal{Q}}$ for each moment that a malicious node stays in the network. A node may be removed by \mathcal{E} at any time, for example through black-listing. However, re-inserting a removed node is not normally possible.

We use N to denote the (possibly random) time at which \mathcal{E} removes a node from the network. In addition, any honest node may leave the network at some (random) time H . Specifically, we assume that an honest node may decide to leave the network with some small probability $\lambda > 0$, independently over time. Then it holds that $\mathbb{E}[H] = \frac{1}{\lambda}$.

Assumption 1: We assume that each node has a fixed type (i.e. honest or malicious) that is not changing over time. The type is hidden from \mathcal{E} .

This does not mean that a node cannot *behave* maliciously during one period and honestly the next: a node that drops packets on purpose, might not do so all the time. Of course, \mathcal{E} not only does not know the type of each node, but it also never observes the rewards obtained or the cost incurred.

¹Although of course the reward is used in the experiments to measure performance.

²Conversely, if we could observe the rewards, it would be trivial to identify malicious nodes.

³This is equivalent to the most likely state approximation and to a sequential probability ratio test under some conditions.

At each time-step t and for each node i , \mathcal{E} receives an *information* signal $x_{i,t} \in [0, 1]$, characterising the behaviour of that node i within the time interval $t \in \mathbb{N}$. This signal can be seen as the output from some IDS, summarising the behaviour of that node during that period.

Assumption 2: We assume that $x_{i,1}, \dots, x_{i,t}$ are independent, (but not identically) distributed, random variables and:

$$\mathbb{E}[x_{i,t} | \mathcal{Q}] = q, \quad \mathbb{E}[x_{i,t} | \mathcal{U}] = u. \quad (1)$$

While the expected value is constant for all t , the observed average of $\frac{1}{t} \sum_{k=1}^t x_{i,k}$ for each node i will initially be far from the expected value for small t . The average, together with the total number of observations for each node form a summary of the information received by each node. The relationship between these quantities will be looked at more closely in the analysis.

Finally, we place no specific meaning to q and u in this work, as they are application-dependent. In an *intrusion response* scenario (e.g. [13]), they could be considered as the *detection rate* (DR) and the *false alarm rate* (FA) correspondingly of an employed intrusion detection system. Then $x_{i,t}$ would correspond to alarm signals, with lower and high values for innocent-looking and suspicious behaviour respectively. Correspondingly, in a *peer-to-peer* scenario (e.g. [17]), they could be fairness or reputation scores of each node.

In the remainder, we always refer to some arbitrary node in the network and thus make no distinction between nodes. This is because the algorithms that we examine, consider each node *independently* of the others. Consequently, the following section analyses the expected loss for a single node of unknown type.

III. THE LOSS MODEL

As previously mentioned, \mathcal{E} obtains a small gain for each time-step an honest node is within the network, and a small loss for each time-step a malicious node remains in the network. Formally, we can write that the total gain G we obtain from some node i , which \mathcal{E} removes at time N , and which would voluntarily leave at time H is:

$$G(i, H, N) = \begin{cases} -N\ell_{\mathcal{Q}}, & i \in \mathcal{Q} \\ \min\{H, N\}g_{\mathcal{U}}, & i \in \mathcal{U}. \end{cases} \quad (2)$$

\mathcal{E} wants to choose some node removal policy π that maximises his total expected gain. That means that \mathcal{E} needs to keep as many as possible honest nodes in the network and eliminate the nodes that behave maliciously. In our analysis, we compare the expected gain of our policy π with that of an *oracle*. The oracle always knows the type of each node (i.e. honest or malicious), and thus, employs the optimal policy π^* . For $i \in \mathcal{Q}$, according to the optimal policy π^* it holds $N = 0$, while for $i \in \mathcal{U}$ according to the optimal policy π^* it is $N = \infty$. Correspondingly,

$$\mathbb{E}_{\pi^*}[G(i)] = \begin{cases} 0, & i \in \mathcal{Q} \\ \mathbb{E}[H]g_{\mathcal{U}}, & i \in \mathcal{U}. \end{cases} \quad (3)$$

Let the loss L be the difference between the gain of the optimal policy and our policy. In particular, the *expected loss* of policy π for a node of type v is defined as:

$$\mathbb{E}_{\pi}[L | v] = \mathbb{E}_{\pi^*(v)}[G | v] - \mathbb{E}_{\pi}[G | v], \quad (4)$$

where the i subscript has been dropped for simplicity. The expected loss is bounded by the *worst-case expected loss*:

$$\mathbb{E}_{\pi}[L] \leq \max_{v \in \{\mathcal{Q}, \mathcal{U}\}} \mathbb{E}_{\pi}[L | v], \quad (5)$$

which we wish to minimise. If \mathcal{E} removes node i from the network at random time N , then he does not receive any more observations $x_{i,t}$ for this node from the IDS. Thus, in essence, we want to find a *stopping rule*, that will let \mathcal{E} to determine the random time N at which stopping occurs, i.e. \mathcal{E} takes the decision that $i \in \mathcal{Q}$ and removes it from the network. We note that, $0 \leq N \leq \infty$, where $N = \infty$ if stopping never occurs.

Since \mathcal{E} does not know if node i is honest or malicious, it must collect a sufficient number of samples so as to only remove nodes for which it is reasonably certain that they are malicious. On the other hand, malicious nodes must be removed as soon as possible, since the operator incurs a cost for every moment they remain in the network. The first algorithm we consider uses simple statistics to make nearly optimal decisions about which nodes to keep.

IV. THE HIPER ALGORITHM

The algorithm, depicted in Alg. 1, uses the knowledge we have about malicious and honest nodes (see equation 1). This is done by calculating the average of all the observations generated by a node i until time t :

$$\theta_t \triangleq \frac{1}{t} \sum_{k=1}^t x_{i,k}, \quad (6)$$

and adding an appropriate *confidence interval* so that errors are made with low probability. Informally, HIPER keeps nodes in the network as long as the statistic θ_t is sufficiently far from the expected statistic q of malicious nodes. In order to avoid throwing away honest nodes prematurely, it always keeps nodes for a certain number of steps to obtain more reliable statistics. However, as time passes, it needs more and more evidence to kick a node out. Consequently, the probability that an honest node is thrown out is bounded.

The analysis of the algorithm proceeds in three steps. First, we calculate the expected loss of the algorithm when faced with a node of malicious type. Then, we calculate the loss for honest nodes. Subsequently, we combine the two losses and tune the algorithm's input parameters to obtain an overall loss bound.

The first bound only depends upon the input parameter δ , the error probability we wish to accept, and the loss $\ell_{\mathcal{Q}}$ incurred by malicious nodes. We prove that the expected loss is polynomially bounded in terms of both δ and $\ell_{\mathcal{Q}}$.

Lemma 1: For Algorithm 1, with input parameter δ , and $\Delta = |u - q|$, the *expected loss* when the node is malicious is

Algorithm 1 HiPER Algorithm for Optimal Response

Parameters: $\delta, \Delta, q \in [0, 1]$

Loop: For each node i in the network:

 For each time-step t do:

 if $|\theta_t - q| < \sqrt{\frac{\ln(2/\delta)}{2t}}$ and $t > \frac{\ln(2/\delta)}{2\Delta^2}$ then

 remove node i from the network

 else keep node i in the network.

 end if

end For

end For

bounded as:

$$\mathbb{E}[L | \mathcal{Q}] \leq \frac{\ell_{\mathcal{Q}}}{(1 - \delta)^2} \quad (7)$$

The proof of this lemma can be found in the appendix. Naturally, the expected loss is linearly dependent on the loss of keeping a malicious node in the network, while the dependence on the error probability is quadratic.

The second bound depends on the input parameter Δ , which corresponds to how far we expect the statistics of honest nodes to be from q , the gain obtained by honest nodes $g_{\mathcal{U}}$ and the leaving probability of honest nodes λ . Once more, we obtain a polynomial loss bound in terms of those variables.

Lemma 2: If $\Delta = |u - q|$, then the *expected loss* when the node is honest is bounded by:

$$\mathbb{E}[L | \mathcal{U}] \leq \frac{g_{\mathcal{U}}(\Delta^2 + 2)}{\lambda(\Delta^2 + 2\lambda)}. \quad (8)$$

The proof of this lemma can be found in the appendix. Similarly to the previous lemma, there is a linear dependence on the loss that is incurred when we erroneously remove an honest node, and a quadratic dependence on the rate of departure. In addition, there is a weak dependence on the gap Δ between the two means.

Finally, we can combine everything in one bound by selecting a value for δ that depends on Δ and which simultaneously makes the bounds tight:

Theorem 1: Set $\Delta = |u - q|$ and select:

$$\delta = 1 - \sqrt{\frac{\ell_{\mathcal{Q}}\lambda(\Delta^2 + 2\lambda)}{g_{\mathcal{U}}(\Delta^2 + 2)}} \quad (9)$$

then the expected loss $\mathbb{E}L$ is bounded by:

$$\mathbb{E}(L) \leq \mathcal{L}_1 \triangleq \frac{g_{\mathcal{U}}(\Delta^2 + 2)}{\lambda(\Delta^2 + 2\lambda)}. \quad (10)$$

Proof: If we substitute (9) in (4) we get:

$$\mathbb{E}[L | \mathcal{Q}] \leq \frac{\ell_{\mathcal{Q}}}{(1 - \delta)^2} = \frac{\ell_{\mathcal{Q}}}{\frac{\ell_{\mathcal{Q}}\lambda(\Delta^2 + 2\lambda)}{g_{\mathcal{U}}(\Delta^2 + 2)}} = \frac{g_{\mathcal{U}}(\Delta^2 + 2)}{\lambda(\Delta^2 + 2\lambda)}. \quad (11)$$

Thus, using (5) and Lemmas 1 and 2 we get:

$$\mathbb{E}(L) \leq \frac{g_{\mathcal{U}}(\Delta^2 + 2)}{\lambda(\Delta^2 + 2\lambda)}$$

This theorem shows that the performance of HiPER only very weakly depends on the gap Δ between honest and malicious nodes. In addition, it is optimal up to a polynomial factor. ■

V. MARKOV DECISION PROCESS APPROXIMATIONS

As mentioned in the introduction, our setting corresponds to an optimal stopping problems. As these can be modelled as MDPs [7], it may be useful to solve the problem directly using the MDP formalism.

To cast our problem in this setting, we need to specify: a) the *prior* probability for each node being *honest* or *malicious*; b) a *known* distribution family for the observation distribution, conditioned on whether the node under consideration is honest or malicious; c) a planning algorithm for calculating our responses. This can be quite demanding computationally, as the solution to the problem requires planning in a large tree. However, they can result in much better performance.

A. Intrusion Response and POMDP

A Partially Observable Markov Decision Process (POMDP) [18] is a generalisation of a Markov Decision Process (MDP). More precisely, a POMDP models the relationship between an agent and its environment when the agent cannot directly observe the underlying state. A POMDP can be described as a tuple $\langle S, A, O, T, \Omega, R \rangle$ where S is a finite set of states, A is a set of possible actions, O is a set of possible observations, T is a set of conditional transition probabilities and Ω is a set of conditional observation probabilities and $R : A, S \rightarrow \mathbb{R}$.

We can model our intrusion response problem as a POMDP if we consider that a node of the network at each time-step t has a state $s_t \in S$ with $s_t = (v_t, c_t)$ where $v_t \in \{0, 1\}$ and $c_t \in \{0, 1\}$ such that:

$$v_t = \begin{cases} 0, & \text{if the node is honest,} \\ 1, & \text{if the node is malicious.} \end{cases}$$

$$c_t = \begin{cases} 0, & \text{if the node is in the network,} \\ 1, & \text{if the node is out of the network.} \end{cases}$$

where it holds that $\mathbb{P}(v_{t+1} = v_t) = 1$ since v_t is stationary (i.e. a malicious node is always malicious and an honest node remains honest) based on Assumption 1.

Additionally, at each time-step t , \mathcal{E} can perform an action $a_t \in \{0, 1\}$ such that:

$$a_t = \begin{cases} 0, & \text{if } \mathcal{E} \text{ keeps the node in the network,} \\ 1, & \text{if } \mathcal{E} \text{ removes the node from the network.} \end{cases}$$

Furthermore, the following independence condition holds: $\mathbb{P}(v_{t+1} | v_t, c_t, a_t) = \mathbb{P}(v_{t+1} | v_t)$ since the type of a node (i.e. malicious or honest) does not depend on \mathcal{E} 's action (i.e. remove from the network or not) neither on whether the node is in the network or out. In addition, since the type of a node never changes, it holds:

$$\mathbb{P}(v_{t+1} = j | v_t = j) = 1. \quad (12)$$

Consequently, we remove the time subscript from v in the sequel. On the other hand the probability that a node will be in the network depends on if it is already in or out and the action that \mathcal{E} will take:

$$\mathbb{P}(c_{t+1} | c_t, v, a_t) = \mathbb{P}(c_{t+1} | c_t, a_t) \quad (13)$$

From equations (12) and (13), it is evident that the POMDP under consideration is factored.

To fully specify the model we must assume some probability distribution for the observations. Specifically, we model x_t as drawn from a Bernoulli distribution⁴ with parameters u and q for honest and malicious nodes respectively: $\mathbb{P}(x_t = 1 | v = 0) = u$ and $\mathbb{P}(x_t = 1 | v = 1) = q$. Let $\mathbf{x}_t \triangleq (x_1, \dots, x_t)$ be a t -length sequence of observations. From Bayes' theorem, we obtain an expression for our *belief* at time t :

$$\mathbb{P}(v = j | \mathbf{x}_t) = \frac{\mathbb{P}(\mathbf{x}_t | v = j) \mathbb{P}(v = j)}{\sum_{i=0}^1 \mathbb{P}(\mathbf{x}_t | v = i) \mathbb{P}(v = i)} \quad (14)$$

where $j \in \{0, 1\}$. Thus, the expected gain at time t if \mathcal{E} decides to keep a node in the network is: $\mathbb{E}[G_t | c_t = 0, \mathbf{x}_t] = \mathbb{P}(v = 0 | \mathbf{x}_t) \cdot g_u - \mathbb{P}(v = 1 | \mathbf{x}_t) \cdot \ell_Q$ while the expected gain if \mathcal{E} decides to remove the node from the network is always: $\mathbb{E}[G_t | c_t = 1] = 0$. The problem is to find a policy $\pi : X^* \rightarrow A$, mapping from the set of all possible sequences of observations to actions, maximising the total expected gain:

$$\mathbb{E}_\pi(G) = \mathbb{E}_\pi \left(\sum_{t=1}^{\infty} G_t \right). \quad (15)$$

Since future gains depend on any future observations we might obtain, the exact calculation requires enumerating all possible future observations. Consequently, the exact solution to the problem is intractable [7], [11], [9]. In the next section we describe possible approximations to this problem.

B. POMDP algorithms

We consider three algorithms: a) A *myopic* algorithm, which only considers the expected gain at the current time-step; b) An *optimistic* algorithm, which computes an upper bound on the total expected gain; c) A *finite lookahead* algorithm, which performs complete planning up to some fixed finite depth. While these algorithms have appeared before in the general MDP literature, they have not been applied before to intrusion response problems. We do not consider the most likely state approximation (MLS), since in our case there are only two possible hidden states for a node, thus, rendering the approximation far too coarse for it to be effective.

1) *Myopic*: In this case, \mathcal{E} only considers the expected gain for the next time-step when taking a decision. Consequently, \mathcal{E} keeps the node in the network if: $\mathbb{E}[G_t | a_t = 0] > \mathbb{E}[G_t | a_t = 1]$. This algorithm is the closest to the MLS approximation among the ones considered. In fact, it is easy to see that it would be identical to MLS, as well as to a sequential probability ratio test, when $\ell_Q = g_u$.

⁴This distribution is particularly convenient for computational reasons, because closed-form Bayesian inference can be performed via the Beta conjugate prior [7]. However, in principle it can be replaced with any other distribution family, without affecting the overall formalism.

2) *Optimistic*: This rule constructs an upper bound on the value of the decision to keep a node in the network, which is based on Proposition 1 in [9]. Informally, this is done by assuming that the true type of the node will be revealed at the next time-step. Then \mathcal{E} keeps the node in the network if and only if: $\mathbb{P}(v_t = 0 | \mathbf{x}_t) \cdot g_u / \lambda > \mathbb{P}(v_t = 1 | \mathbf{x}_t) \cdot \ell_Q$. Intuitively, if the node is revealed to be malicious, then we can remove it at the next step and consequently we only lose ℓ_Q . In the converse case, we can keep it for an expected $1/\lambda$ steps.

3) *Finite lookahead*: The finite lookahead algorithm performs backwards induction [7] up to some finite depth T , at every time-step. More precisely, any sequence of observations $\mathbf{x}_t = (x_1, \dots, x_t)$ results in a posterior probability $\mathbb{P}(v_t | \mathbf{x}_t)$. Let: $V_t \triangleq \sum_{k=t}^{\infty} G_k$ be the total gain starting from time-step t . Then, the expected gain under the optimal policy is determined recursively as follows:

$$\begin{aligned} \mathbb{E}(V_t | \mathbf{x}_t) &= \max\{0, \mathbb{E}(G_t | \mathbf{x}_t, a_t = 0) + \mathbb{E}(V_{t+1} | \mathbf{x}_t)\} \\ \mathbb{E}(V_{t+1} | \mathbf{x}_t) &= p_t \mathbb{E}(G_t | \mathbf{x}_t, x_{t+1} = 1) + \\ &\quad (1 - p_t) \mathbb{E}(V_{t+1} | \mathbf{x}_t, x_{t+1} = 0) \end{aligned}$$

where $p_t \triangleq \mathbb{P}(x_{t+1} = 1 | \mathbf{x}_t) = \sum_{i=0}^1 \mathbb{P}(x_{t+1} = 1 | v = i) \mathbb{P}(v = i | \mathbf{x}_t)$ is the marginal posterior probability that $x_{t+1} = 1$. For more details on this backwards induction algorithm, the reader is urged to consult [7], [11].

VI. EXPERIMENTAL EVALUATION

We perform three sets of experiments. The first set investigates the performance of HiPER with various choices of the parameter δ , including the optimal choice suggested by Theorem 1. The second set compares HiPER with the *myopic* and *optimistic* approximations. In the final set of experiments, we compare the *optimistic* with the *finite lookahead* approximation. In all cases, we collected results from 10^4 runs, with 100 nodes in each simulation, and we plot a moving average of the *expected loss* as various network parameters change. Specifically, the first results we report (i.e. Fig. 1) are made through 10^4 experiments. For each experiment, we selected a horizon $H \sim \text{Uniform}([10, 1000])$, user and adversary parameters $u, q \sim \text{Uniform}([0, 1])$, and user gain $g_u \sim \text{Uniform}([0, 1])$ and we set $\ell_Q = 1$. Each experiment measured the loss for a network containing 100 nodes, each of which had a probability p of being malicious, with $p \sim \text{Beta}(2, 2)$ for each experiment. During each run, the i -th node generates a sequence of observations $x_{i,t}$ drawn from a Bernoulli distribution with parameter u if the node is honest and q if the node is malicious. Figure 1 shows a summary of the results, averaged over these trials. It can be seen that, while HiPER's performance is relatively robust to the choice of δ , nevertheless the optimal choice suggested by Theorem 1 generally leads to small losses.

For our second set of experiments, shown in Figure 2, we compare HiPER with the *optimistic* and *myopic* algorithms. We increased the range of user gains to $g_u \sim \text{Uniform}([0, 2])$ compared to the previous setup, but the other experimental parameters remain the same. It is clear that the *myopic*

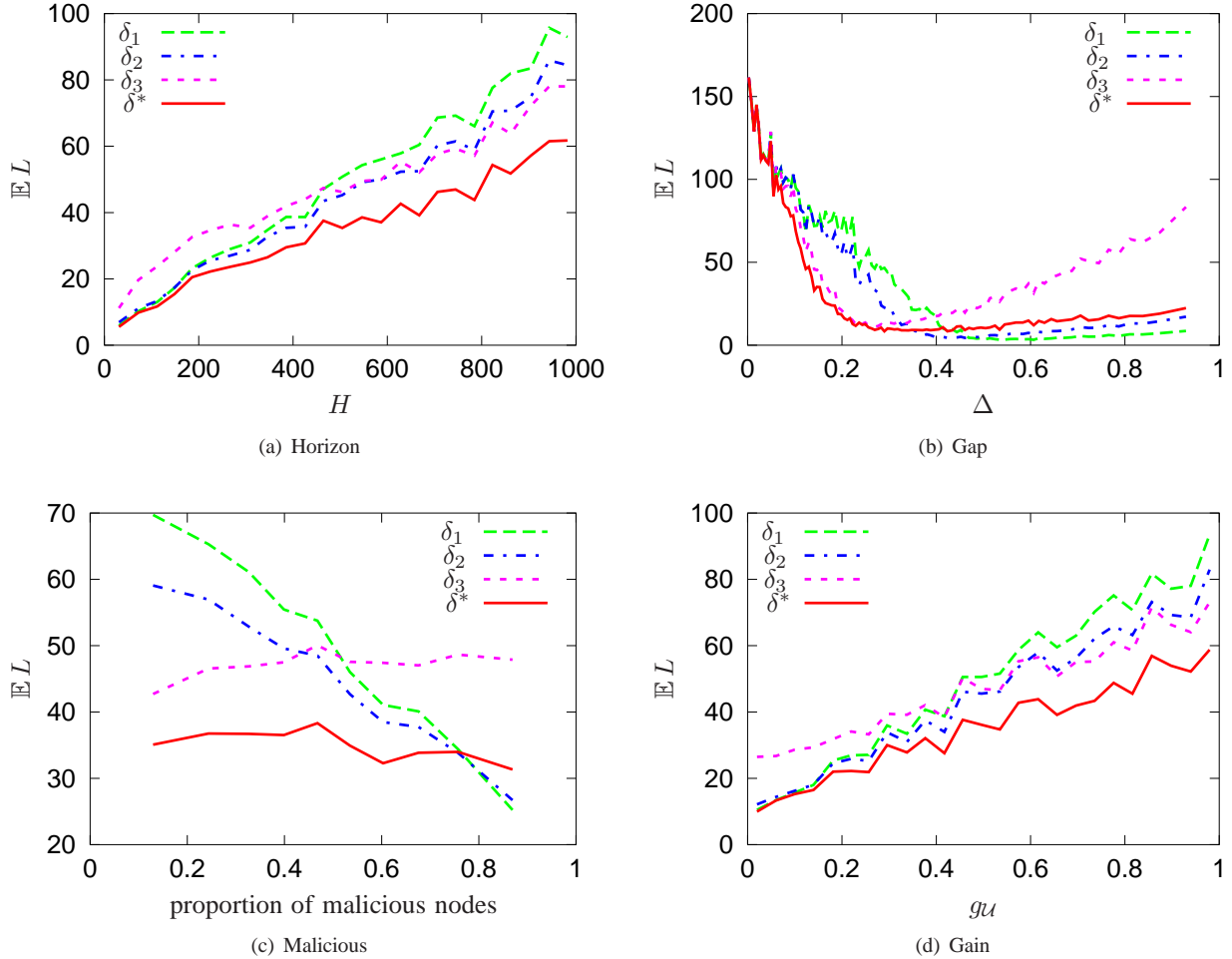


Fig. 1. Simulations with Alg. 1, for four different choices of δ . In particular $\delta_1 = 0.9$, $\delta_2 = 0.95$, $\delta_3 = 0.99$ and δ^* is chosen according to Theorem 1. It can be seen that, while the algorithm is not extremely sensitive to the exact choice of δ , the optimal value is generally more robust.

approximation has almost always a higher loss compared to both HiPER and the *optimistic* algorithm. The latter, while performing at a similar level to HiPER, has an advantage when either the proportion of malicious is small or when gu is large. This makes sense intuitively, since in those cases the optimism is justified. In the converse case, however, the *optimistic* approach performs worse than HiPER, which is less sensitive to the proportion of malicious nodes, since it is a worst-case approach.

Finally, we performed some experiments comparing the *optimistic* approximation with the *finite-lookahead* POMDP solvers for lookahead for T time-steps where $T \in \{4, 8\}$. While these do not solve the problem to the end of the horizon H , they plan ahead for T steps at every time-step of the simulation. Unfortunately, the complexity of these solvers is exponential in T , which limited the amount of simulations we could perform to 10^3 and we only considered horizons $H \sim \text{Uniform}([1, 100])$. These experiments are shown in Fig. 3. In comparison with Fig. 2, the *finite lookahead* algorithms performs much better than the *myopic* approximation and indeed the 8-step lookahead manages to

slightly outperform the *optimistic* approximation. In addition, it is much more robust to the proportion of malicious nodes in the network. However, the relative advantage of the 8-step to the 4-step lookahead is relatively small for the amount of extra computation required.⁵

VII. CONCLUSION

This paper defined a resource management problem that arises frequently in communication networks. Namely, whether to remove a suspicious node from the system, with the amount of available evidence, or to collect some further data before taking the final decision. This is in fact a type of stopping problem, which we believe is of relevance to many applications where blacklisting may be performed. This includes applications such as automated intrusion response, as well as ensuring fairness in peer-to-peer networks, such as [19]. To this end, we proposed and analysed, both theoretically and experimentally, an efficient algorithm, HiPER, that achieves low *worst-case expected loss* relative to an oracle that knows *a priori* the type (honest or malicious) of every

⁵The computational effort is exponential in T .

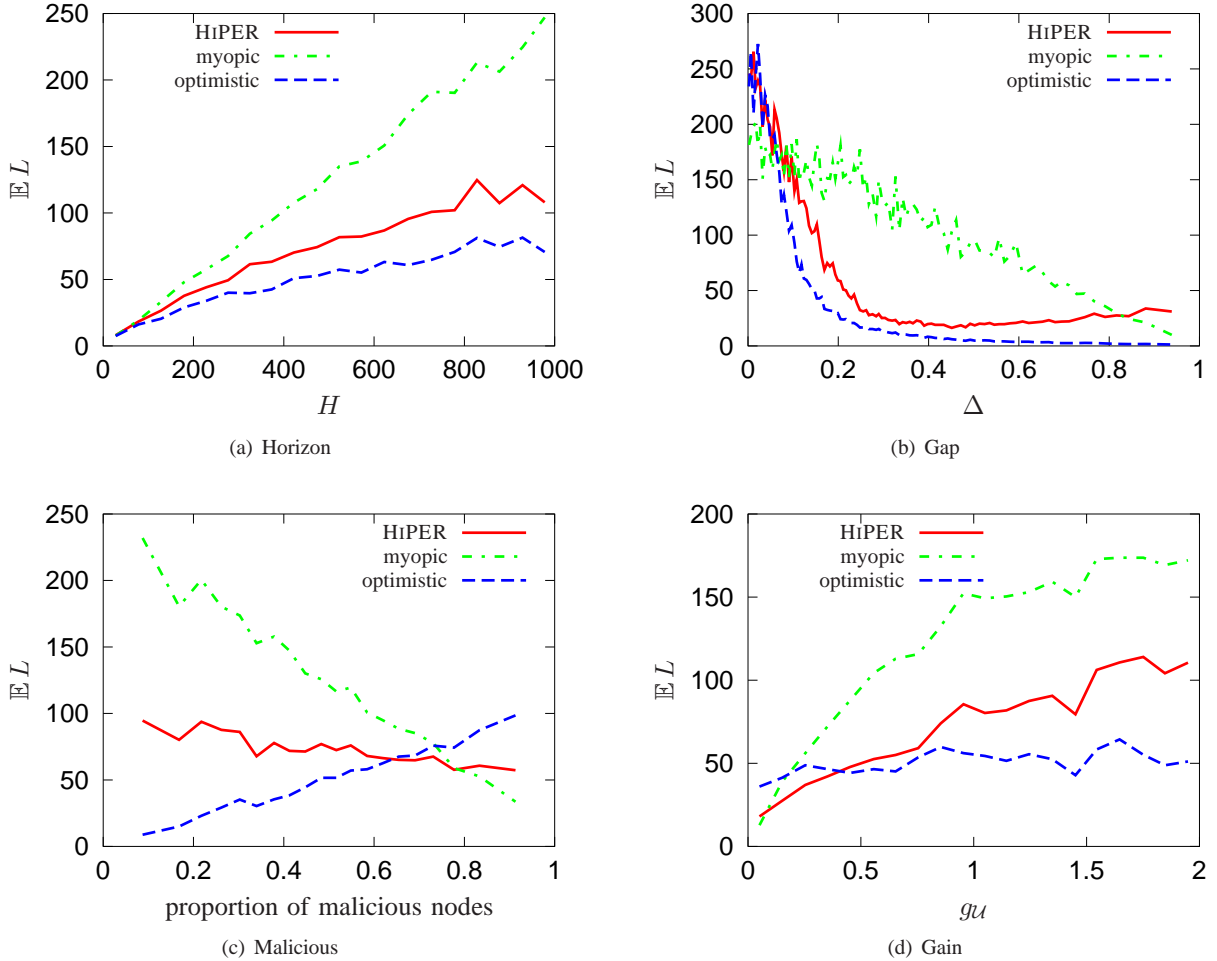


Fig. 2. Comparison of HiPER with the *myopic* solver and the *optimistic* approximation for various network conditions. It can be clearly seen that the *myopic* approximation is significantly worse than both approaches. However, the *optimistic* approach outperforms the worst-case HiPER algorithm when the proportion of malicious nodes is low. The *optimistic* approach is also better when the payment for honest nodes is high.

node in the system. In addition, we derived and compared a number of algorithms by modelling the problem as a POMDP: a *myopic* and an *optimistic* approximation, as well as a *finite lookahead* solver. Of those, the *optimistic* approximation and the partial *finite lookahead* solvers perform the best, with the *finite lookahead* methods being the most robust, while simultaneously being computationally demanding.

The main advantage of HiPER are its simplicity and lack of stringent assumptions on the distribution. This makes it suitable for deployment in most situations. However, whenever a full probabilistic model and computational resources are available, one of the approximate solvers would be useful. The overall best performance is offered by the *finite lookahead*, closely followed by the *optimistic* approximation. The *myopic* approximation, which is equivalent to the widely-used “most likely state” (MLS) approximation, is the worst. To our knowledge, neither the *optimistic* approximation, nor the *finite lookahead* methods have been applied before to this problem or more generally to intrusion response problems. They should be more generally applicable for other types

of intrusion response and resource management problems. It is our view that they are inherently more suitable than other approximations such as the commonly used (MLS) approximation (or equivalently, a sequential probability ratio test) which in our setting produces an essentially random policy.

For future work, we would like to extend our theoretical analysis to the performance of the *optimistic* and the *finite lookahead* algorithms. In addition, it would be interesting to examine a more general game-theoretic scenarios, including strategic attackers [1], [10]. Finally, we would like to generalise our setting so that observations must be *explicitly* gathered from each node, where it is not possible to continuously sample all nodes due to budget constraints. In fact, the sampling problem in the context of intrusion detection, has been recently studied by [14], [3]. A natural extension of our work would consequently be to optimally combine sampling and response policies.

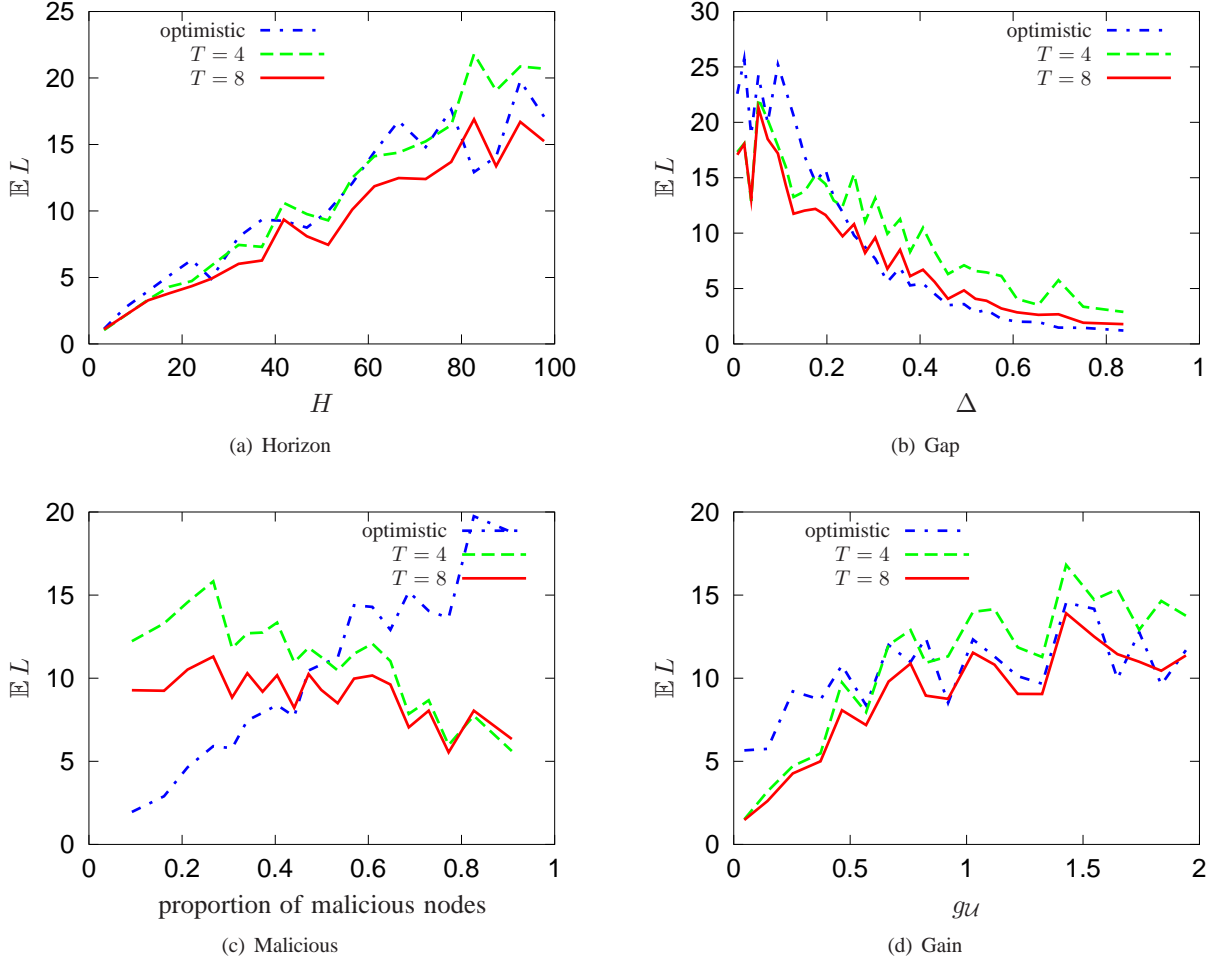


Fig. 3. Comparison of the optimistic approximation with approximate *non-myopic* POMDP solvers for planning lookahead of T time-steps where $T \in \{4, 8\}$. It can be seen that, for short horizons, these perform just as well and that they are more robust to the proportion of malicious nodes in the network. However, these methods are computationally more intensive, with complexity $O(e^T)$.

APPENDIX

This section collects the missing proofs from the main text.

(*Lemma 1*): Since the node i under consideration is malicious, i.e. $i \in \mathcal{Q}$, it holds that: $\mathbb{E}[x_{i,t} | \mathcal{Q}] = q$. Then, we have:

$$\mathbb{E}[\theta_t | \mathcal{Q}] = \mathbb{E}\left[\frac{1}{t} \cdot \sum_{k=1}^t x_{i,k} \mid \mathcal{Q}\right] = \frac{1}{t} \sum_{k=1}^t \mathbb{E}[x_{i,k} | \mathcal{Q}] = \frac{1}{t} \cdot t \cdot q = q. \mathbb{E}[x_{i,t} | \mathcal{U}] = u.$$

From Hoeffding's inequality (*Lemma 3*, in the Appendix), we have:

$$\mathbb{P}(|\theta_t - q| > \epsilon_t | \mathcal{Q}) \leq 2 \exp(-2t\epsilon_t^2), \quad (16)$$

where $\epsilon_t > 0$ and $\mathbb{P}(|\theta_t - q| > \epsilon_t | \mathcal{Q})$ denotes the probability that θ_t (which is random) is very far away from q (which is fixed). Now set: $\epsilon_t = \sqrt{\frac{\ln(2/\delta)}{2t}}$ as in Algorithm 1. Then, since equation 16 holds for any $\epsilon_t > 0$, we get that the probability of keeping a malicious node $i \in \mathcal{Q}$ in the network is at most δ : $\mathbb{P}\left(|\theta_t - q| > \sqrt{\frac{\ln(2/\delta)}{2t}} \mid \mathcal{Q}\right) < \delta$. Thus, we have: $\mathbb{E}[L |$

$$\mathcal{Q}] = \mathbb{E}[N | \mathcal{Q}] \cdot \ell_{\mathcal{Q}} = \sum_{t=0}^{\infty} \mathbb{P}(N = t | \mathcal{Q}) \cdot t \cdot \ell_{\mathcal{Q}} \leq \ell_{\mathcal{Q}} \sum_{t=0}^{\infty} \delta^{t-1} \cdot t = \frac{\ell_{\mathcal{Q}}}{(1-\delta)^2} \quad \blacksquare$$

(*Proof of Lemma 2*): We denote by N the time-step at which \mathcal{E} removes node i from the network. Then, the function $g : \mathbb{N}^2 \rightarrow \mathbb{R}$ that gives us the gain for each node i is defined as: $g(n, h) \triangleq \min\{n, h\} \cdot g_{\mathcal{U}}$ where $h \in H$ and $n \in N$. Since the node i under consideration is honest, i.e. $i \in \mathcal{U}$, we have $u = q + \Delta$, where $\Delta > 0$. So we only need $\mathbb{P}(\theta_t - q < \epsilon_t | \mathcal{U})$. Since $q = u - \Delta$ from the Hoeffding inequality (*Lemma 3*, in the Appendix), we have:

$$\begin{aligned} \mathbb{P}(N = t | \mathcal{U}) &\leq \mathbb{P}(N \leq t) \leq \mathbb{P}(\theta_t - u < \epsilon_t - \Delta | \mathcal{U}) \\ &\leq \exp(-2 \cdot t(\epsilon_t - \Delta)^2) \end{aligned}$$

where $\Delta - \epsilon_t > 0$. It holds that:

$$\begin{aligned} \mathbb{E}[G | \mathcal{U}, N = n] &= \\ \sum_{n=0}^{\infty} \mathbb{P}(H = h | \mathcal{U}, N = n) \mathbb{E}[G | \mathcal{U}, N = n, H = h] \quad (17) \end{aligned}$$

But it holds that: $\mathbb{E}[G \mid \mathcal{U}, N = n, H] = g(n, h)$ and since $h \in H$ and $n \in N$ are independent we have: $\mathbb{P}(H = h \mid \mathcal{U}, N = n) = \mathbb{P}(H = h \mid \mathcal{U})$. Thus,

$$\begin{aligned} \mathbb{E}[G \mid \mathcal{U}, N = n] &= \sum_{h=0}^{\infty} \mathbb{P}(H = h \mid \mathcal{U}) \cdot g(n, h) = \\ &= \sum_{h=0}^{\infty} \mathbb{P}(H = h \mid \mathcal{U}) \min\{n, h\} \cdot g_{\mathcal{U}} = \\ &= g_{\mathcal{U}} \cdot \left\{ \sum_{h=0}^{n-1} \mathbb{P}(H = h \mid \mathcal{U}) \cdot h + \sum_{h=n}^{\infty} \mathbb{P}(H = h \mid \mathcal{U}) \cdot n \right\} \quad (18) \end{aligned}$$

The expected loss is given by subtracting from the expected gain of the oracle policy, when \mathcal{E} never removes the node from the network (i.e. $N = \infty$), the expected gain when \mathcal{E} removes the node at the time-step $N = n$. Thus, it holds:

$$\begin{aligned} \mathbb{E}[L \mid \mathcal{U}, N = n] &= \mathbb{E}[G \mid \mathcal{U}, N = \infty] - \mathbb{E}[G \mid \mathcal{U}, N = n] = \\ &= \lim_{n \rightarrow \infty} (\mathbb{E}[G \mid \mathcal{U}, N = n]) - \mathbb{E}[G \mid \mathcal{U}, N = n] = \\ &= g_{\mathcal{U}} \sum_{h=0}^{\infty} \mathbb{P}(H = h)h - g_{\mathcal{U}} \left\{ \sum_{h=0}^{n-1} \mathbb{P}(H = h)h + \sum_{h=n}^{\infty} \mathbb{P}(H = h)n \right\} \\ &= g_{\mathcal{U}} \left\{ \sum_{h=n}^{\infty} \mathbb{P}(H = h)h - \sum_{h=n}^{\infty} \mathbb{P}(H = h)n \right\} \quad (19) \end{aligned}$$

Since, by definition $\mathbb{P}(H = h + 1 \mid H > h) = \lambda$, we have $\mathbb{P}(H = h) = (1 - \lambda)^{h-1} \lambda$. Consequently,

$$\begin{aligned} \mathbb{E}[L \mid \mathcal{U}, N = n] &= g_{\mathcal{U}} \lambda \left(\sum_{h=n}^{\infty} (1 - \lambda)^{h-1} h - \sum_{h=n}^{\infty} (1 - \lambda)^{h-1} n \right) \\ &= g_{\mathcal{U}} \frac{(1 - \lambda)^n}{\lambda} \quad (20) \end{aligned}$$

Thus, we have: $\mathbb{E}[L \mid \mathcal{U}] = \sum_{t=0}^{\infty} \mathbb{P}(N = t \mid \mathcal{U}) \mathbb{E}[L \mid N = t] \leq \sum_{t=0}^{\infty} \exp(-2 \cdot t \cdot (\epsilon_t - \Delta)^2) \cdot g_{\mathcal{U}} \frac{(1 - \lambda)^t}{\lambda}$. Since the algorithm uses $\epsilon_t = \frac{\Delta}{\sqrt{t}}$, we have:

$$\begin{aligned} \mathbb{E}[L \mid \mathcal{U}] &\leq \frac{g_{\mathcal{U}}}{\lambda} \sum_{t=0}^{\infty} \exp\left(-2 \cdot t \left[\frac{\Delta}{\sqrt{t}} - \Delta\right]^2\right) (1 - \lambda)^t \\ &= \frac{g_{\mathcal{U}}}{\lambda} \sum_{t=0}^{\infty} \exp\left(-2\Delta^2(\sqrt{t} - 1)^2\right) (1 - \lambda)^t \\ &\leq \frac{g_{\mathcal{U}}}{\lambda} \sum_{t=0}^{\infty} \exp\left(-2\Delta^2 \left[\sqrt{t} - \sqrt{\frac{t}{2}}\right]^2\right) (1 - \lambda)^t \\ &= \frac{g_{\mathcal{U}}}{\lambda} \sum_{t=0}^{\infty} \left[\exp\left(-\frac{\Delta^2}{2}\right) (1 - \lambda)\right]^t \\ &= \frac{g_{\mathcal{U}}}{\left[1 - \exp\left(-\frac{\Delta^2}{2}\right) (1 - \lambda)\right] \lambda} \leq \frac{g_{\mathcal{U}}(\Delta^2 + 2)}{\lambda(\Delta^2 + 2\lambda)} \end{aligned}$$

where $t \geq 2$. \blacksquare

Definition 1 (Bernoulli distribution): If X_1, \dots, X_n are independent Bernoulli random variables with $X_k \in \{0, 1\}$ and $\mathbb{P}(X_k = 1) = \mu$ for all k , then

$$\mathbb{P}\left(\sum_{k=1}^n X_k \geq u\right) = \sum_{k=0}^u \binom{n}{k} \mu^k (1 - \mu)^{n-k}. \quad (21)$$

Lemma 3 (Hoeffding): For independent random variables X_1, \dots, X_n such that $X_i \in [a_i, b_i]$, with $\mu_i \triangleq \mathbb{E} X_i$ and $t > 0$:

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq \sum_{i=1}^n \mu_i + nt\right) \leq \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

The same inequality holds for $\sum_{i=1}^n X_i \leq \sum_{i=1}^n \mu_i - nt$.

REFERENCES

- [1] N. Bao, P. Kreidl, and J. Musacchio. A network security classification game. In *GameNets 2011*, 2011.
- [2] C. Boutilier. A POMDP formulation of preference elicitation problems. In *Proceedings of the National Conference on Artificial Intelligence*, pages 239–246. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2002.
- [3] S. Bu, F. Yu, X. Liu, and H. Tang. Structural results for combined continuous user authentication and intrusion detection in high security mobile ad-hoc networks. *Wireless Communications, IEEE Transactions on*, (99):1–10, 2011.
- [4] T. Bui, M. Poel, A. Nijholt, and J. Zwiers. A tractable DDN-POMDP approach to affective dialogue modeling for general probabilistic frame-based dialogue systems. 2006.
- [5] A. Cassandra. *Exact and Approximate Algorithms for Partially Observable Markov Decision Processes*. PhD thesis, Brown University, 1998.
- [6] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning and Games*. 2006.
- [7] M. DeGroot. *Optimal Statistical Decisions*. John Wiley & Sons, 1970. Republished in 2004.
- [8] S. Dejmal, A. Fern, and T. Nguyen. Reinforcement learning for vulnerability assessment in peer-to-peer networks. In *Proceedings of the 20th national conference on Innovative applications of artificial intelligence*, pages 1655–1662, 2008.
- [9] C. Dimitrakakis. Complexity of stochastic branch and bound methods for belief tree search in Bayesian reinforcement learning. In *2nd international conference on agents and artificial intelligence (ICAART 2010)*, pages 259–264, Valencia, Spain, 2009. ISNTICC, Springer.
- [10] L. Dritsoula, P. Loiseau, and J. Musacchio. A game-theoretical approach for finding optimal strategies in an intruder classification game. In *CDC 2012*, 2012.
- [11] M. O. Duff. *Optimal Learning Computational Procedures for Bayes-adaptive Markov Decision Processes*. PhD thesis, University of Massachusetts at Amherst, 2002.
- [12] Y. He and K. Chong. Sensor scheduling for target tracking in sensor networks. In *Decision and Control, 2004. CDC. 43rd IEEE Conference on*, volume 1, pages 743–748. IEEE, 2004.
- [13] W. Lee, W. Fan, M. Millerand, S. Stolfo, and E. Zadok. Toward Cost-Sensitive Modeling for Intrusion Detection and Response. *Journal of computer Security*, 10:5–22, 2000.
- [14] K. Liu and Q. Zhao. Dynamic intrusion detection in resource-constrained cyber networks. Technical Report arXiv:112.0101, arXiv, 2011.
- [15] S. Ross, J. Pineau, S. Paquet, and B. Chaib-draa. Online planning algorithms for POMDPs. *Journal of Artificial Intelligence Resesarch*, 32:663–704, July 2008.
- [16] Z. Saigol and U. of Birmingham. School of Computer Science. *Information-lookahead planning for AUV mapping*. School of Computer Science, University of Birmingham, 2009.
- [17] P. Si, F. Yu, H. Ji, and V. Leung. Distributed sender scheduling for multimedia transmission in wireless mobile peer-to-peer networks. *Wireless Communications, IEEE Transactions on*, 8(9):4594–4603, 2009.
- [18] R. Smallwood and E. Sondik. The Optimal Control of Partially Observable Markov Processes over a Finite Horizon. *Operational Research*, 21:1071–88, 1973.
- [19] A. Vieira, S. Campos, and J. Almeida. Fighting attacks in P2P live streaming: Simpler is better. In *INFOCOM Workshops 2009, IEEE*, pages 1–2. IEEE, 2009.
- [20] X. Zan, F. Gao, J. Han, X. Liu, and J. Zhou. A Hierarchical and Factored POMDP based Automated Intrusion Response Framework. In *Proceedings of the 2nd International Conference on Software Technology and Engineering (ICSTE)*, volume 2, pages 410–414. IEEE, 2010.
- [21] Z. Zhang, P.-H. Ho, and L. He. Measuring IDS-estimated Attack Impacts for Rational Incident Response: A Decision Theoretic Approach. *Computers & Security*, 28:605–614, 2009.

- [22] S. Zonouz, H. Khurana, W. Sanders, and Y. T.M. RRE: A Game-Theoretic Intrusion Response and Recovery Engine. In *Proceedings of the IEEE/IFIP International Conference on Dependable Systems & Networks, 2009 (DSN'09)*, pages 439–448, Lisbon, Portugal, 29 June – 2 July 2009.