

Adversarial Online Learning with Variable Plays in the Pursuit-Evasion Game: Theoretical Foundations and Application in Connected and Automated Vehicle Cybersecurity

Yiyang Wang *Graduate Student Member, IEEE*, Neda Masoud

Abstract—We extend the adversarial/non-stochastic multi-play multi-armed bandit (MPMAB) to the case where the number of arms to play is variable. The work is motivated by the fact that the resources allocated to scan different critical locations in an interconnected transportation system change dynamically over time and depending on the environment. By modeling the malicious hacker and the intrusion monitoring system as the attacker and the defender, respectively, we formulate the problem for the two players as a sequential pursuit-evasion game. We derive the condition under which a Nash equilibrium of the strategic game exists. For the defender side, we provide an exponential-weighted based algorithm with sublinear pseudo-regret. We further extend our model to heterogeneous rewards for both players, and obtain lower and upper bounds on the average reward for the attacker. We provide numerical experiments to demonstrate the effectiveness of a variable-arm play.

Index Terms—adversarial bandit, cyber security, pursuit-evasion game, online learning, intelligent transportation systems (ITS), multi-armed bandit (MAB), algorithmic learning theory

I. INTRODUCTION

Currently, the world is experiencing an evolution from the traditional transportation system to the next generation of intelligent transportation systems (ITS). ITS aims to satisfy the ever-increasing need for mobility in major cities, which has caused growing traffic congestion, air pollution, poor user experience and crashes. Developing a sustainable intelligent transportation system requires better usage of the existing infrastructures and their seamless integration with information and communication technologies (ICT). Enabled by the recent findings in the areas of telecommunications, electronics, and computing capabilities in recent decades, the subsystems (infrastructures and vehicles) in ITS are expected to interoperate and communicate with each other, in order to provide a better and safer traveling experience [1].

The interconnection between the infrastructures and the vehicles relies on various types of sensors to provide state information and situational awareness. However, this has also increased the vulnerability of these advanced systems to cyber attacks. For instance, recently there have been demonstrated cyber attacks on vehicle sensors in [2], [3], where the authors used optimization-based approaches to fool the light detection

and Ranging (LiDAR) sensors on the vehicle. At the system level, the infrastructures and the vehicles can be viewed as individual nodes in a large interconnected network, where a single malicious attack on a subset of sensors of one node can easily propagate through this network, affecting other network components (e.g., other vehicles, traffic control devices, etc.). For example, Feng et al. [4] demonstrated that by sending falsified data to actuated and adaptive signal control systems, a malicious hacker could increase the total system delay in a real-world corridor. Therefore, there is an increasing need for cyber security solutions, especially for sensor security solutions, to enhance the safety and reliability of the entire system.

Cyber security is an extremely broad topic. However, previous work on cyber security in the realm of ITS mainly focuses on either attack or the defense strategies. For instance, there exists a large body of research illustrating the potential risks of connected and automated vehicle (CAV) technologies that result in anomalous/false information [5]–[8]. In the case of CAV sensor security, several critical sensors are illustrated in [9], including differential global positioning systems (GPS), inertial measurement units, engine control sensors, tyre-pressure monitoring systems (TPMS), LiDAR, and camera. Meanwhile, CAVs require more engine control units (ECUs) and many features of CAVs require complex interactions between multiple ECUs, which may potentially expose more vulnerabilities compared to non-CAVs. There also exist several studies assessing the potential threats on the transportation infrastructure [4], [10], [11]. For example, field devices such as traffic signals and roadside units are susceptible to tampering. The aforementioned literature illustrates the potential threats of sensor attacks to connected transportation systems.

Besides threat detection, prevention is normally recognized as one of the best defense strategies against malicious hackers or attackers. In order to deploy better prevention mechanisms, behaviors of both the attacker and the defender have to be considered so that the attack profile can be predicted. There is a gap in the literature in considering both the attacker and the defender and the adaptive interactions between them when devising defense strategies, which this paper aims to bridge.

Moreover, as more sensors are mounted aboard CAVs or installed on the transportation infrastructure, it becomes more difficult to monitor the sensors continuously, mainly due to limited resources. Although there is a large body of literature addressing sensor security in ITS [12]–[16], most of them mainly focus on sensor intrusion/anomaly detection without attack profile analysis, which considers which sensor is more

vulnerable and should be protected. In this study, we address this by modeling attacker and defender behaviors in a game theoretical framework. Specifically, instead of considering intrusion/anomaly detection for all sensors in the system, we model attack and defense behaviors in order to predict which subset of sensors are more likely to be compromised. To be more practical, we consider a dynamic resource constraint for the defender. We model this problem as a sequential evasion-and-pursuit game between two players. Consider the intrusion monitoring system of a sensor network as the defender. At each time, the defender selects a subset of sensors to scan, while the number of selected sensors changes based on the environment and scanning history, among other factors. Meanwhile, a hacker, considered as the attacker, attempts to select a sensor to compromise without being scanned by the defender. We assume that both the attacker and the defender are able to learn their opponent's behavior adaptively and with only partial information over time, and investigate the resulting decision problem.

The main contributions of this work are as follows: First, in order to predict the attack profile, we model the behaviors of the attacker and the defender as the adversarial (or non-stochastic) multi-armed bandit (MAB) problem and the multi-armed bandit problem with variable plays (MAB-VP), where the two players are playing a constant-sum game against each other. To the best of our knowledge, this is the first study of MAB-VP in the non-stochastic setting. Second, we derive conditions under which a Nash equilibrium of the strategic game exists. For the defender, we provide an exponential-weighted algorithm, which is shown to have sublinear pseudo-regret. Finally, we consider a more realistic setting where the rewards are heterogeneous among different sensors, and derive lower and upper bounds on the attacker's average reward.

II. LITERATURE REVIEW

In this paper, we explore online learning algorithms in the class of adversarial or non-stochastic multi-armed bandit (MAB) problems. The adversarial MAB problem was first addressed by Auer et al. [17], where they also proposed the well-known exponential-weight algorithm for exploration and exploitation (Exp3). Exp3 runs the Hedge algorithm, which was originally proposed by Freund and Schapire [18] as a subroutine. Since then, there have been several extensions to this class including the online shortest path problem [19], routing games [20], bandit online linear optimization [21], and combinatorial bandits [22].

The multi-play multi-armed bandit (MPMAB) problem is another research direction for MAB. In this extension, a fixed number of resources (i.e., arms) are allocated at each time step. The MPMAB has attracted a lot of interest and several studies have been conducted along this direction [23]–[26]. However, most of these studies only focus on a stochastic setting. There is much less concentration on the adversarial MPMAB problem: Cesa-Bianchi and Lugosi [22] considered combinatorial bandits in the adversarial setting, where they proposed the ComBand algorithm. This algorithm has a sublinear regret in $O\left(M^{\frac{3}{2}}N\sqrt{TN\ln N}\right)$, with time and space complexities of

$O(MN^3)$ and $O(K^3)$, respectively, where M is the number of resources (or arms selected) at each time, T is the number of iterations, and N is the number of possible actions. Following this work, Uchiya et al. [27] proposed an extension of Exp3, Exp3.M, which runs in $O(N(\log M + 1))$ time and $O(N)$ space, and suffers at most $O\left(\sqrt{MTN\log(N/M)}\right)$ regret. However, the aforementioned algorithms only consider a fixed number of arms to be played at each time.

Only a limited number of studies have considered variable plays. Fouché et al. [28] proposed a scaling algorithm combined with a MAB algorithm, which they call the S-MAB algorithm. In this algorithm, the number of arms played at each time changes in order to satisfy an efficiency constraint. However, although the authors considered a dynamic environment, the S-MAB algorithm uses a stochastic setting, where they assume an unknown distribution of reward for each arm. Another work addressing the variable plays problem was done by Lesage-Landry and Taylor [29], where they extended the stochastic MAB to stochastic plays setting, i.e. the number of arms to play evolves as a stationary process. Both these studies only considered a stochastic setting, and did not conduct any game strategy analysis.

Although there is a wealth of research on using game theory in the transportation literature, very few studies applied game theory in ITS cybersecurity. Sedjelmaci et al. [30] conducted a survey on recent studies utilizing game theory to protect ITS from attacks, which is to the best of our knowledge the only survey paper on this topic. However, without considering the adaptive behavior of opponents, the current literature mostly models the cybersecurity problem as a non-repeated game, such as the Stackelberg security games (SSG) [31], [32], zero-sum games [33], [34], or Bayesian games [35], [36]. The solutions from these types of models are typically in the form of equilibria with an implied assumption that the players have knowledge of their opponent's actions/beliefs. Instead, we formulate this cybersecurity problem as a sequential pursuit-evasion game, which is also in the realm of algorithmic learning theory. There have been several studies of the pursuit-evasion problem [37]–[39]. However, they either lack robustness against adaptive changes in the adversarial behavior, or do not consider multiple plays, variable plays, dynamic resource allocation, or heterogeneous rewards.

Since the behavior of the adversarial opponent usually cannot be described in a stochastic way, in this paper we study the MAB-VP problem in a non-stochastic setting, where we propose the Exp3.M with variable plays (Exp3.M-VP) algorithm. Next, we consider a game setting for two players, and show that a Nash equilibrium of the strategic game exists. Finally, we consider heterogeneous rewards for both players and derive lower and upper bounds for the attacker's average reward. Numerical analyses are conducted in order to further demonstrate our results.

III. SYSTEM MODEL AND PROBLEM FORMULATION

A. System Model

Consider the repeated pursuit-evasion game between an attacker and a defender in discrete time. At each time step

TABLE I: Table of Notation

$\alpha_k(t)/\beta_k(t)$	\triangleq	marginal probability that the attacker compromises/the defender scans location k at time t
$x_k(t)/y_k(t)$	\triangleq	indicator variable of whether the defender/attacker selects the location k at time t
I_t/J_t	\triangleq	index of the locations where the attacker compromises/the defender scans at time t
M_t	\triangleq	number of locations scanned by the defender at time t
a/b	\triangleq	lower/upper bound of M_t
$r(t)/s(t)$	\triangleq	single step reward of the attacker/defender
$\omega(t)/\theta(t)$	\triangleq	private randomization device of the attacker/defender
π_t/γ_t	\triangleq	control policy of the attacker/defender
T	\triangleq	finite time horizon
N	\triangleq	total number of locations
\mathcal{N}	\triangleq	index set of N locations
\mathcal{C}	\triangleq	index set of arbitrary locations

t , the attacker selects one of the N locations, indexed by the set $\mathcal{N} = \{1, 2, \dots, N\}$, to hide in (e.g., compromise a sensor), while the defender searches M_t locations simultaneously, where $1 \leq a \leq M_t \leq b < N$. The behaviors of the attacker and the defender are described by their respective set of marginal probabilities $\alpha(t) = (\alpha_k(t))_{k \in \mathcal{N}}$ and $\beta(t) = (\beta_k(t))_{k \in \mathcal{N}}$, where $\alpha_k(t)$ and $\beta_k(t)$ are the respective probabilities that the k -th location is chosen by the attacker and the defender at time t . Note that $\alpha(t)$ and $\beta(t)$ represent the adversarial behavior with respect to one's opponent at time t , where they can describe randomized strategies of the players, or a probabilistic belief held by one side about the likelihood of an action by the other side.

Define two sets of binary variables $x_k(t)$ and $y_k(t)$ such that $x_k(t) = 1$ if the defender does not search location k at time t , and $x_k(t) = 0$ otherwise. Similarly, $y_k(t) = 1$ if the attacker compromises the location k at time t , and $y_k(t) = 0$ otherwise. When the attacker (defender) does not know the type of algorithm/strategy the opponent uses, it may regard the $x_k(t)$ ($y_k(t)$) as a predetermined but unknown number. When the attacker (defender) does have this information, it may regard the $x_k(t)$ ($y_k(t)$) as a random variable, where $P(x_k(t) = 0) = \beta_k(t)$ (resp. $P(y_k(t) = 1) = \alpha_k(t)$). The game is played in a sequence of trials $t = 1, 2, \dots, T$. In this work we consider the case that neither the attacker nor the defender knows the strategy adopted by the other player. As will be discussed later, they have to choose the location based on the their history rewards.

B. Problem Formulation: Partial Information Game

In this study we consider the scenario where both players have limited information on the adaptive behavior of their opponent. Define $\pi = (\pi_t, t = 1, 2, \dots)$ as the control policy of the attacker, and let Π denote the policy space. Denote the location selection (action) sequence as $I = (I_t, t = 1, 2, \dots)$ under policy π and $|I_t| = 1$. At each time and under policy π_t , the attacker chooses one location $I_t \in \mathcal{N}$ to attack, i.e.,

$$I_t = \pi_t \left(x_I^{[t-1]}, I^{[t-1]}, \omega(t) \right), \quad (1)$$

where $x_I^{[t-1]} := (x_{I_1}(1), \dots, x_{I_{t-1}}(t-1))$, and $I^{[t-1]}$ is similarly defined. $(\omega(t), t = 1, 2, \dots)$ denotes the randomized

strategy of the attacker. Let $x_k(t)$ be the state of location k for the attacker at time t . Then the attacker scores the corresponding reward $r^I(t) = x_{I_t}(t)$. The attacker observes only the reward $r^I(t)$ for the chosen action I_t .

The attacker receives an expected reward $E[r^I(t)] = 1 - \beta_{I_t}(t)$ at time t , which is the mean number of successful attacks at the chosen location. Note that in this section we consider a homogeneous reward across all locations; however, heterogeneous location-dependent rewards are considered in section VII. In this study, we assume a 100% success rate for both attacks and detection attempts. Then, within the time window $\{t, t = 1, 2, \dots, T\}$, the attacker considers the following maximization problem,

$$\underset{\pi \in \Pi, I_t \in \mathcal{N}}{\text{maximize}} \mathbb{E} \left\{ \sum_{t=1}^T x_{I_t}(t) \right\}, \quad (2)$$

where the expectation is with respect to the randomness of the system state and the mixed-strategy of the attacker.

We assume that the defender can scan M_t locations at time t . Define $\gamma = (\gamma_t, t = 1, 2, \dots)$ as the control policy of the defender, and let Γ denote the defender's policy space. Denote the location selection (action) sequence as $J = (J_t, t = 1, 2, \dots)$ under policy γ . At each time and under policy γ_t , the defender scans M_t locations, denoted as set $J_t \subset \mathcal{N}$ and $|J_t| = M_t$, based on their history search and rewards, i.e.,

$$J_t = \gamma_t \left(y_J^{[t-1]}, J^{[t-1]}, M_t, \theta(t) \right), \quad (3)$$

where $y_J^{[t-1]} := (y_{J_1}(1), \dots, y_{J_{t-1}}(t-1))$ with $J^{[t-1]}$ similarly defined, and $(\theta(t), t = 1, 2, \dots)$ denotes the randomized strategy of the defender. Let $y_k(t)$ be the state of location k for the defender at time t . The defender also observes only the rewards $\sum_{j \in J_t} y_j(t)$ of the selected action J_t . Denote the total rewards at time t of the defender given the location selection sequence J as $s^J(t) = \sum_{j \in J_t} y_j(t)$. The defender therefore receives the expected reward $E[s^J(t)] = \sum_{j \in J_t} \alpha_j(t)$ at time t . This expected reward represents the mean number of detected attacks among M_t number of scanned locations.

We assume that the number of arms M_t the defender plays at each time is determined by a scaling function, i.e. $f : \mathbb{R}^{N+1} \rightarrow \{a, a+1, \dots, b\}$, of the d -moving average of the rewards of each arm, where a and b are integers,

and $1 \leq a \leq b < N$. We also assume that M_t is a function of the environment constraint L_t , since in reality checking a location (e.g., scanning a specific sensor/unit in a CAV) may consume resources. Then, given the time horizon T , the defender is trying to solve the following constrained optimization problem:

$$\underset{\gamma \in \Gamma, J_t \subset \mathcal{N}}{\text{maximize}} \mathbb{E} \left\{ \sum_{t=1}^T \sum_{j \in J_t} y_j(t) \right\} \quad (4a)$$

$$\text{s.t. } M_t = f(\hat{y}^d(t), L_t) \quad (4b)$$

$$|J_t| = M_t \quad (4c)$$

where $\hat{y}^d(t) := (\hat{y}_1^d(t), \hat{y}_2^d(t), \dots, \hat{y}_N^d(t))$, and \hat{y}_i^d is the d -moving average of the rewards of each arm i . Using a moving average of reward can allow us to capture the history reward while at the mean time capturing the dynamic change of the reward for each location, allowing the scaling function to adjust the number of arms to play each time. The expectation is with respect to the randomness of the system state and the mixed strategy of the defender. Note that there is no requirement for the scaling function f , other than it needs to be bounded by integers a and b . Furthermore, L_t can be an arbitrary integer between a and b , thereby capturing any set of environmental conditions.

When the defender knows the type of strategy the attacker uses, it may regard $y_j^J(t)$ as stochastic, i.e. assuming the attacker chooses location j with probability $P(y_j^J = 1) = \alpha_j(t)$. Note that this is different from the stochastic MAB setting where a fixed (time-invariant) distribution of rewards for each arm is assumed. However, here we do not assume neither the defender nor the attacker have information about their opponent's strategy. Hence, the difficulty is that the defender can only estimate $\alpha_j(t)$ by imposing an arbitrary belief on the adversarial behavior based on previous observations and rewards. Furthermore, here, we do not make any assumptions about the distribution of $\alpha_j(t)$.

IV. ALGORITHMS FOR THE ATTACKER AND THE DEFENDER

We assume the attacker adopts Exp3 proposed by Auer et al [17]. (However, as we are going to show later in section VI, the equilibrium of the two-player game does not depend on any properties of the algorithm other than a no-regret guarantee.) The Exp3 algorithm uses an efficient and randomized policy to select only one arm at each time t . The adversarial single play bandit problem is closely related to the problem of learning to play an unknown repeated matrix game. In this setting, a player without prior knowledge of the game matrix is to play the game repeatedly against an adversary with complete knowledge of the game and unbounded computational power. The basic idea of Exp3 is that at each time the player uses a randomized policy such that the adversarial player cannot know the exact choice of the player before she/he plays. For the details of Exp3, refer to the Appendix A.

Unlike the attacker who selects a single location to attack, we assume the defender can search multiple number of locations, which may vary at each time. Both sides seek to

Algorithm 1 Exp3.M-VP

```

1: Parameter:  $\eta \in (0, 1]$ 
2: Initialization:  $w_i(1) = 1$  for  $i = 1, 2, \dots, N$ 
3: for  $t = 1, 2, \dots, T$  do
4:   Receive the number of arms to play at each round  $M_t$ .
5:   if  $\max_{j \in \mathcal{N}} w_j(t) \geq \left(\frac{1}{M_t} - \frac{\eta}{N}\right) \sum_{i=1}^N w_i(t)/(1-\eta)$  then
6:     Decide  $\kappa_t$  such that
       
$$\frac{\kappa_t}{\sum_{w_i(t) \geq \kappa_t} \kappa_t + \sum_{w_i(t) < \kappa_t} w_i(t)} = \left(\frac{1}{M_t} - \frac{\eta}{N}\right)/(1-\eta).$$

       Set  $S_0(t) = \{i : w_i(t) \geq \kappa_t\}$ .
       Set  $w'_i(t) = \kappa_t, \forall i \in S_0(t)$ .
7:   else
8:     Set  $S_0(t) = \emptyset$ .
9:   end if
10:  Set  $w'_i(t) = w_i(t), \forall i \in S_0^c(t)$ .
11:  Set  $\hat{\alpha}_i(t) = M_t \left( (1-\eta) \frac{w'_i(t)}{\sum_{j=1}^N w'_j(t)} + \frac{\eta}{N} \right)$ .
12:  Set  $J_t = \text{DepRound}(M_t, (\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_N))$ .
13:  Observe rewards  $y_i(t) \in [0, 1]$  for  $i \in J_t$ .
14:  for  $i = 1, 2, \dots, N$  do
15:    
$$\hat{y}_i(t) = \begin{cases} y_i(t)/\hat{\alpha}_i(t) & \text{if } i \in J_t, \\ 0 & \text{otherwise.} \end{cases}$$

    
$$w_i(t+1) = \begin{cases} w_i(t) \exp(M_t \eta \hat{y}_i(t)/N) & \text{if } i \in S_0^c(t), \\ w_i(t) & \text{otherwise.} \end{cases}$$

16:  end for
17: end for
```

maximize their respective total rewards. At the beginning of a time step, each side needs to decide which location(s) to target, and cannot change their selection until the next time step. We develop a variable-play extension of the Exp3.M algorithm for the defender, which we call Exp3.M-VP, as detailed in Algorithm 1. In the Exp3.M-VP algorithm, let S denote the set of selected locations, and let S^c define its complement set. Under the non-stochastic assumption and at each time step, the Exp3.M-VP algorithm consists of the following two procedures:

- 1) Receive M_t , which is determined by the scaling function f and could be based on the environment constraint L_t as well as the history rewards $\hat{y}^d(t)$ at time t , among other factors. Note that function f can take any form, and defining its exact form is outside the scope of this paper. Here, we assume M_t is provided.
- 2) Apply an adversarial MPMAB algorithm which selects M_t arms (locations) to play.

For the second procedure, we use the Exp3.M algorithm as a subroutine of the Exp3.M-VP algorithm. The Exp3.M is proposed by Uchiya et al. [27] and is an extension of the algorithm Exp3 for the adversarial MPMAB setting. In contrast to the Exp3 algorithm which selects one arm at each time, Exp3.M randomly selects a fixed number of M arms at each time. Note that both Exp3 and Exp3.M suffer from sublinear (weak) regret, or no-regret. In order to make sure that the probability of selecting location i by DepRound at step 12, i.e. $\hat{\alpha}_i(t)$, does not exceed 1, the Exp3.M-VP algorithm checks whether all $w_j(t)$'s are less than $\left(\frac{1}{M_t} - \frac{\eta}{N}\right) \sum_{i=1}^N \frac{w_i(t)}{(1-\eta)}$ at

step 5. If that is the case, $\hat{\alpha}_i(t)$ calculated at step 11 will be less than 1 for all $i = 1, 2, \dots, N$ without any weight modification, and the set $S_0(t)$ is set to \emptyset at step 8. Otherwise, all the actions i with $w_i(t) \geq \kappa_t$ are classified into $S_0(t)$ and set to κ_t at step 6. Doing this, we have $\hat{\alpha}_i(t) = 1$ for all $i \in S_0(t)$. The subroutine **DepRound** [40] at step 12 draws M_t out of N items with the specified marginal distribution $(\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_N)$, and is included in Appendix B.

V. ADAPTIVE LEARNING OF THE DEFENDER

In this section, we address the adaptive learning of the defender. Based on Algorithm 1 for the defender, the problem (4) can be recast by removing the constraint set, since will divide the problem to a scaling procedure and the MAB-VP. Formally, let $\mathbf{y}(t) := (y_k(t), \forall k \in \mathcal{N})$ for $t = 1, \dots, T$ over a finite horizon T . For any search sequence of the defender $J = (J_t, t = 1, 2, \dots)$ and a fixed sequence of attacks by the attacker $(\mathbf{y}(1), \mathbf{y}(2), \dots)$, the total reward of the defender at T , denoted by $G^J(T)$, is given by

$$G^J(T) = \sum_{t=1}^T \sum_{j \in J_t} y_j(t). \quad (5)$$

Here, we obtain the maximum reward by consistently searching the subset \mathcal{A}_{M_t} , which is the most attacker-active location set at each time step t with cardinality M_t :

$$G_{\max}(T) = \max_{\mathcal{A}_{M_t}} \sum_{t=1}^T \sum_{k \in \mathcal{A}_{M_t}} y_k(t). \quad (6)$$

Let us define $\mathcal{A} = \cup_{M_t} \mathcal{A}_{M_t}$, where $\mathcal{A} \subset \mathcal{N}$. Note that if $M_t \in \{a, a+1, \dots, b\}$, the location index subset \mathcal{A}_{M_t} is defined such that $\mathcal{A}_a \subset \mathcal{A}_{a+1} \subset \dots \subset \mathcal{A}_b = \mathcal{A}$.

The regret is then defined as

$$R(T) = G_{\max}(T) - G^J(T). \quad (7)$$

When $a = b$, i.e. M_t is time-invariant, the above regret reduces to the standard regret of MPMAB problem.

Since we care more about the competition against the optimal action in expectation, we define the pseudo-regret for our MAB-VP problem following the definition of pseudo-regret in [41] as:

$$\bar{R}(T) = G_{\max}(T) - E[G^J(T)], \quad (8)$$

where the expectation is with respect to the randomness of the system state and the mixed-strategy of the defender.

Theorem 1. For any $N > 0$ and for any $\eta \in (0, 1]$, if M_t is lower bounded and upper bounded by two positive integers a and b respectively, then

$$\begin{aligned} \bar{R}_{\text{Exp3.M-VP}}(T) &= G_{\max}(T) - E[G_{\text{Exp3.M-VP}}^J(T)] \\ &\leq \left(1 + \frac{(e-2)b}{a}\right) \eta G_{\max}(T) + \frac{N}{\eta} \ln \frac{N}{b} \end{aligned} \quad (9)$$

holds for any assignment of rewards and for any $T > 0$.

Proof. See Appendix C. ■

By appropriately choosing the parameter η , we can obtain the following corollary:

Corollary 1.1. Set $\eta = \min \left\{ 1, \sqrt{\frac{Na \ln(N/b)}{(a+(e-2)b)bT}} \right\}$. Then

$$\bar{R}_{\text{Exp3.M-VP}}(T) \leq 2 \sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \sqrt{bTN \ln \frac{N}{b}}$$

holds for any $T > 0$ and for any assignment of rewards.

The proof of Corollary 1.1 is the same as that of Corollary 3.2 in [42]. For the proof of Corollary 1.1, see Appendix D. Note that when $a = b$, the upper bound in Corollary 1.1 is the same as the upper bound of Exp3.M in [27], and when $a = b = 1$ the upper bound becomes the same upper bound obtained for Exp3 in [17].

Corollary 1.2. Define $\bar{s}_{\infty} := \liminf_{T \rightarrow \infty} E \left[\frac{1}{T} \sum_{t=1}^T s^J(t) \right]$ as the average reward of the defender over infinite time horizon. Using the same parameter η as in Corollary 1.1, when the defender uses the Exp3.M-VP algorithm against the attacker who adopts a no-regret algorithm, we have $\bar{s}_{\infty} = \frac{\nu}{N}$ if M_t is a wide sense stationary process with mean ν .

In order to prove Corollary 1.2, we need the following lemma, which was originally derived in [39].

Lemma 2. When the defender (pursuer) is adopting Exp3.M and the attacker (evader) does not know the type of algorithm used by the adversarial opponent, then $v = \frac{1}{N}$, where v is the game value of the repeated constant-sum game for the defender.

Then the proof of Corollary 1.2 is as follows.

Proof. The above problem is equivalent to the problem of two players playing an unknown repeated bimatrix game, where the game value $v_{i,t}$ ($i = 1, 2$ for the row and column player respectively) is changing over time. Define the game matrices as two $N \times N$ matrices \mathbf{B} and \mathbf{C} , where $B_{ij} + C_{ij} = 1$ for any $(i, j) \in \mathcal{N} \times \mathcal{N}$. At each time t , the defender (i.e., the row player) chooses J_t rows of the matrix, and at the same time, the attacker (i.e., the column player) chooses exactly one column $I_t = k$. The defender then receives the payoff $\sum_{j \in J_t} B_{jk} = \sum_{j \in J_t} y_j(t)$. The defender uses a mixed strategy \mathbf{p}_t at each time t , where $\mathbf{p}_t \in [0, 1]^N$, and the attacker chooses according to a probability vector $\mathbf{q}_t \in [0, 1]^N$. Note that the sum of \mathbf{p}_t equals M_t and the sum of \mathbf{q}_t equals 1. Let $v_{1,t}$ be the game value of the game matrix \mathbf{B} at time t . Then by Corollary 1.1, we have

$$E \left[\sum_{t=1}^T \sum_{j \in J_t} B_{jk} \right] = E \left[\sum_{t=1}^T \sum_{j \in J_t} y_j(t) \right] \quad (10a)$$

$$\begin{aligned} &\geq G_{\max}(T) - 2 \sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \\ &\quad \times \sqrt{bTN \ln \frac{N}{b}}. \end{aligned} \quad (10b)$$

Let \mathbf{p}_t be such that

$$v_{1,t} = \max_{\mathbf{p}_t} \min_{\mathbf{q}_t} \mathbf{p}_t^\top \mathbf{B} \mathbf{q}_t = \min_{\mathbf{q}_t} \max_{\mathbf{p}_t} \mathbf{p}_t^\top \mathbf{B} \mathbf{q}_t.$$

Then we have

$$G_{\max}(T) \geq \sum_{t=1}^T \sum_{i=1}^N p_{t,i} y_i(t) \quad (11a)$$

$$= \sum_{t=1}^T \mathbf{p}_t^\top \mathbf{y}(t) \quad (11b)$$

$$= \sum_{t=1}^T \mathbf{p}_t^\top \mathbf{B} \mathbf{q}_t \geq \sum_{t=1}^T v_{1,t} \quad (11c)$$

where \mathbf{q}_t is a distribution vector whose I_t -th component is 1.

Combining (10) and (11), we have

$$\begin{aligned} E \left[\frac{1}{T} \sum_{t=1}^T s^J(t) \right] &\geq \frac{1}{T} \sum_{t=1}^T v_{1,t} - 2 \sqrt{\left(1 + (e-2) \frac{b}{a}\right)} \\ &\quad \times \sqrt{bN \ln \frac{N}{b} / T}. \end{aligned} \quad (12)$$

Note that at each time t , $v_{1,t} = M_t v_1$, where v_1 is the game value when the defender only chooses one location. Hence, by taking the limit of (12) and according to the law of large numbers we have

$$\bar{s}_\infty = \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T v_{1,t} = \nu v_1, \quad (13)$$

where the first equality comes from the fact that the attacker is also adopting a no-regret algorithm (e.g. Exp3). Finally, according to Lemma 2, we obtain the result. \blacksquare

Corollary 2.1. *Under the setting that the defender adopts Exp3.M-VP and the attacker adopts a no-regret algorithm, assuming that M_t is a wide sense stationary process with mean ν , each player adopts the best response for the infinite-horizon problem.*

The proof can be obtained by extending the proof of the defender side in Corollary 1.2 to both sides, and is omitted for brevity. Note that in Corollary 1.2 and Corollary 2.1 we do not specify which type of learning algorithm the attacker is using, and the only assumption is that the attacker adopts a no-regret algorithm.

VI. ADAPTIVE LEARNING OF THE ATTACKER

We assume that the attacker adopts the Exp3 algorithm to randomly attack one location at each time step. The Exp3 algorithm runs the algorithm Hedge as a subroutine. Unlike the Hedge algorithm which directly takes advantage of the full information of the reward vector $\mathbf{x}(t) := (x_i(t), \forall i \in \mathcal{N})$, Exp3 observes partial information and feeds the simulated reward vector $\hat{\mathbf{x}}(t) := (\hat{x}_i(t), \forall i \in \mathcal{N})$ to the Hedge. The Hedge will then update $\hat{\beta}_i(t)$, which is the prediction of probability $\beta_i(t)$ for $i \in \mathcal{N}$. For more details about the Exp3 and Hedge algorithms, see Appendix A.

The defender adopts the Exp3.M-VP algorithm, which has a sublinear regret, as shown in Theorem 1. As a result, if the attacker favors one location, intuitively the defender will eventually identify this most attractive location, and fails to scan it only at a rate no more than sublinear in T . When M_t is a time-invariant constant, it follows immediately that the best strategy for the attacker over an infinite time horizon is to treat each location equally, either in a stochastic or deterministic way. However, when M_t is a variable, the same argument cannot be trivially made.

Theorem 3. *Define $\bar{r}_\infty := \liminf_{T \rightarrow \infty} E \left[\frac{1}{T} \sum_{t=1}^T r^I(t) \right]$, and let the location sequence g be the sequence of the greedy policy π_{greedy} , where $g(t) = \arg \min_{i \in \mathcal{N}} \hat{\beta}_i(t)$ for all t . If M_t is bounded by two positive integers a, b such that $M_t \in \{a, a+1, \dots, b\}$, then under any policy π we have:*

$$\bar{r}_\infty \leq \frac{N-a}{N},$$

and under the greedy policy π_{greedy} ,

$$\bar{r}_\infty \geq \frac{N-b}{N}.$$

Proof. See Appendix E. \blacksquare

Note that by Corollary 1.2, we can directly obtain the following result,

Corollary 3.1. *Under the setting that the defender adopts Exp3.M-VP, the attacker adopts Exp3, and M_t is a wide sense stationary process with mean ν , we have $\bar{r}_\infty = \frac{N-\nu}{N}$.*

Moreover, when M_t is a wide sense stationary process, following the proof of Theorem 3, it is not hard to show that even the greedy policy can obtain $\bar{r}_\infty = \frac{N-\nu}{N}$. Note that the above argument does not require Exp3.M-VP to have any property other than a no-regret guarantee, and therefore the greedy policy for the attacker can be a countermeasure against the entire family of no-regret algorithms. For the defender part, according to Corollary 1.2 and Corollary 3.1, a straightforward path to increase the average reward in an infinite time horizon is to increase the value of ν , i.e., assign more resources to the intrusion monitoring system.

VII. ADAPTIVE ADVERSARIAL LEARNING WITH HETEROGENEOUS REWARDS

In this section we consider heterogeneous rewards that are location-dependent. This corresponds to a more general setting, since in reality some locations (e.g., sensors) are more critical to the system than others. Let μ_k be the location-dependent reward corresponding to the k -th location. That is, the rewards of the attacker and the defender are $r^I(t) = \mu_{I_t} x_{I_t}(t)$ and $s^J(t) = \sum_{j \in \mathcal{J}_t} \mu_j y_j(t)$, respectively. Without loss of generality, we assume that $\mu_1 \geq \mu_2 \geq \dots \geq \mu_N$. We denote the frequency of location k being selected given the selection sequence I as $d_k^I(T)$ over a time horizon T , i.e.,

$$\frac{1}{T} \sum_{t=1}^T \mu_{I_t} = \frac{1}{T} \sum_{k=1}^N c_k^I(T) \mu_k = \sum_{k=1}^N d_k^I(T) \mu_k \quad (14)$$

where $c_k^I(T) = |\{t \leq T : I_t = k\}|$ and $d_k^I(T) = c_k^I(T)/T$. Note that $c_k^I(T)$ is the total number of times location k is selected by the attacker over horizon T given the selection sequence I .

Since the problem is no longer a constant-sum game under the setting of heterogeneous rewards, Corollary 2.1 and Corollary 3.1 cannot be directly applied. However, we can still show that when the reward for each location is heterogeneous, the average reward \bar{r}_∞ in an infinite time horizon is bounded within an interval determined by a , b , and μ_k , $k = 1, 2, \dots, N$.

Theorem 4. *Given heterogeneous rewards, the average reward of the attacker \bar{r}_∞ over an infinite time horizon is bounded within the interval $\left[\frac{K^*-b}{\sum_{k=1}^{K^*} \mu_k}, \frac{K^*-a}{\sum_{k=1}^{K^*} \mu_k}\right]$, where K^* is a constant determined by μ_k values such that $b \leq K^* \leq N$.*

In order to prove Theorem 4, we need Lemmas 5 and 6, as follows. Let $\text{supp}(\mathbf{d}) = \{k \in \mathcal{N} : d_k > 0\}$ for any feasible solution \mathbf{d} , and let K^* be the cardinality of $\text{supp}(\mathbf{d})$. Then we have the following lemmas:

Lemma 5. *For any optimal solution \mathbf{d}^* of problem (18), (i) $\mu_k d_k^* = \mu_j d_j^*$ for any $k, j \in \text{supp}(\mathbf{d}^*)$, and (ii), $\text{supp}(\mathbf{d}^*)$ consists of the indices of locations with the K^* highest μ .*

Lemma 6. *Problem (23) is lower bounded by $\frac{K^*-b}{\sum_{k=1}^{K^*} \mu_k}$.*

The proofs of Lemmas 5 and 6 can be found in the Appendices F and G, respectively.

Now we shall give the proof of Theorem 4 as follows.

Proof. The average reward of the defender when using Exp3.M-VP is given by

$$E[G_{\text{Exp3.M-VP}}^J(T)] = E\left[\sum_{t=1}^T \sum_{j \in J_t} \mu_j y_j(t)\right] \quad (15a)$$

$$= \sum_{t=1}^T \sum_{k=1}^N \mu_k y_k(t) \beta_k(t) \quad (15b)$$

$$= \sum_{t=1}^T \mu_{I_t} \beta_{I_t}(t) \quad (15c)$$

$$= \sum_{t=1}^T \mu_{I_t} - E\left[\sum_{t=1}^T r^I(t)\right] \quad (15d)$$

for any realization I .

Then we have

$$\frac{1}{T} E\left[\sum_{t=1}^T r^I(t)\right] = \frac{1}{T} \sum_{t=1}^T \mu_{I_t} - \frac{1}{T} E[G_{\text{Exp3.M-VP}}^J(T)] \quad (16a)$$

$$\leq \sum_{k=1}^N \mu_k d_k^I(T) - \frac{1}{T} \left(G_{\max}(T) - 2\sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \times \sqrt{bTN \ln \frac{N}{b}} \right) \quad (16b)$$

$$\leq \sum_{k=1}^N \mu_k d_k^I(T) - \max_{J \in \mathcal{C}(N,a)} \sum_{j \in J} \mu_j d_j^I(T) + 2\sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \sqrt{bN \ln \frac{N}{b}}/T \quad (16c)$$

where $\mathcal{C}(\mathcal{N}, a) = \{S \subseteq \mathcal{N} : |S| = a\}$, namely, the set of all subsets of size a in \mathcal{N} . The second inequality uses the fact that

$$\begin{aligned} G_{\max}(T) &\geq \max_{J \in \mathcal{C}(N,a)} \sum_{t=1}^T \sum_{j \in J} \mu_j y_j(t) \\ &= \max_{J \in \mathcal{C}(N,a)} \sum_{j \in J} \mu_j c_j^I(T). \end{aligned}$$

Therefore, by having T approach infinity, we have

$$\bar{r}_\infty \leq \liminf_{T \rightarrow \infty} E\left[\sum_{k=1}^N \mu_k d_k^I(T) - \max_{J \in \mathcal{C}(N,a)} \sum_{j \in J} \mu_j d_j^I(T)\right] \quad (17)$$

for any policy π .

Consider the following optimization problem

$$\text{maximize}_{\mathbf{d} \in \Delta_N} \sum_{k=1}^N \mu_k d_k - \max_{J \in \mathcal{C}(N,a)} \sum_{j \in J} \mu_j d_j, \quad (18)$$

where Δ_N is the set of distributions over \mathcal{N} and $\mathbf{d} = (d_k, k \in \mathcal{N})$. Let the optimal solution and its objective function value be \mathbf{d}^* and r_{\max} , respectively. Then we have

$$\bar{r}_\infty \leq r_{\max} = \sum_{k=1}^N \mu_k d_k^* - \max_{J \in \mathcal{C}(N,a)} \sum_{j \in J} \mu_j d_j^*. \quad (19)$$

Without loss of generality, we assume that $\text{supp}(\mathbf{d}^*) = \{1, 2, \dots, K^*\}$. Therefore, according to Lemma 5, we have $d_k^* = \frac{1/\mu_k}{\sum_{j=1}^{K^*} 1/\mu_j}$ for all $k \leq K^*$. Then the optimal value of problem (18) is given by $(K^* - a)/\sum_{j=1}^{K^*} 1/\mu_j$, which is increasing with respect to the value of $K^* = 1, 2, \dots, N$. This gives the upper bound of \bar{r}_∞ .

When the defender adopts Exp3M-VP, we have

$$E[G_{\text{Exp3}}^I(T)] \geq G'_{\max}(T) - o(T). \quad (20)$$

where $G_{\text{Exp3}}^I(T)$ is the total reward of the attacker when adopting Exp3, and $G'_{\text{max}}(T) = \max_{k \in N} \sum_{t=1}^T x_k(t)$ is the maximum total reward the attacker can gain when selecting a fixed location to attack.

Similarly, define $h_k^J(T) = |\{t \leq T : k \in J_t\}|$ and $l_k^J(T) = h_k^J(T)/T$. Then, we have

$$G'_{\text{max}}(T) = \max_{k \in N} \mu_k(T - h_k^J(T)). \quad (21)$$

Thus, the average reward \bar{r}_∞ of the attacker over an infinite time horizon is lower bounded by

$$\bar{r}_\infty \geq \liminf_{T \rightarrow \infty} E \left\{ \max_{k \in N} \mu_k(1 - l_k^J(T)) \right\}. \quad (22)$$

Consider the following optimization problem

$$\underset{c \in \Delta_N}{\text{minimize}} \max_{k \in N} \mu_k(1 - l_k), \quad (23)$$

and denote the optimal value of problem (23) as r_{\min} . Then according to Lemma 6, $r_{\min} = \frac{K^* - b}{\sum_{k=1}^{K^*} \mu_k}$, which gives us the lower bound of \bar{r}_∞ .

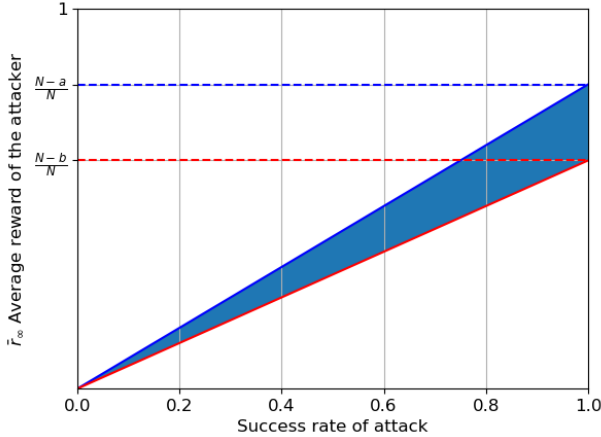


Fig. 1: Range of the average reward of the attacker in an infinite time horizon under different attack success rates.

Theorem 4 is practical when the attack success rate of the attacker is not 100 percent for all locations, where μ_k represents the success rate of attacks on location k for the attacker. Note that although Theorem 4 assumes heterogeneous rewards, it can be simply applied to homogeneous rewards as well. Figure 1 shows the range for the attacker's average reward in an infinite time horizon under different attack success rates, where we assume the same attack success rate for all locations for simpler visualization. Note that we do not even assume that M_t is a wide sense stationary process; the only assumption here is that it is confined within a range with lower and upper bounds a and b , respectively. The shaded blue region in Figure 1 indicates the potential reward the attacker can obtain in infinite time, and the red and blue lines indicate the lower and upper bounds on the attacker's average reward in

infinite time, according to Theorem 4. When the attack success rate is 1, the lower and upper bounds become equivalent to the bounds in Theorem 3. It is straightforward to see that the lower the success rate of the attack, the safer the system will be.

VIII. NUMERICAL ANALYSIS

We conducted extensive simulations illustrating the performance of the proposed algorithm and policy. Our numerical analysis consists of three parts. In section VIII-A, we conduct simulations to test the Exp3.M-VP performance under a single-player setting. In section VIII-B, we compare the performance of Exp3.M-VP with several bandit learning algorithms, i.e., the Exp3, Exp3.M, upper-confidence Bound (UCB) [43], and ϵ -greedy algorithms [44], on real in-vehicle network datasets from the Car-Hacking datasets [45]. In section VIII-C, we run simulations on the proposed game model and algorithmic solutions.

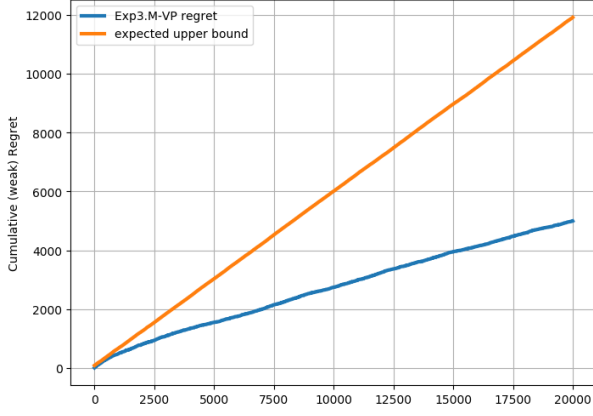
A. Simulations on a Single Player

In this section we consider the single-player setting, where the Exp3.M-VP algorithm was evaluated on a ten-armed bandit problem with rewards for arms drawn independently from Bernoulli distributions with means $\{0.75, \dots, \frac{3}{4k}, \dots, 0.075\}$, with $k = 1, 2, \dots, 10$. This scenario was simulated over a fixed time horizon $T = 20,000$ time steps. The number of arms played at each time step is drawn independently from a discrete uniform distribution over $\{1, 2, 3\}$. Parameter η is set to 0.1.

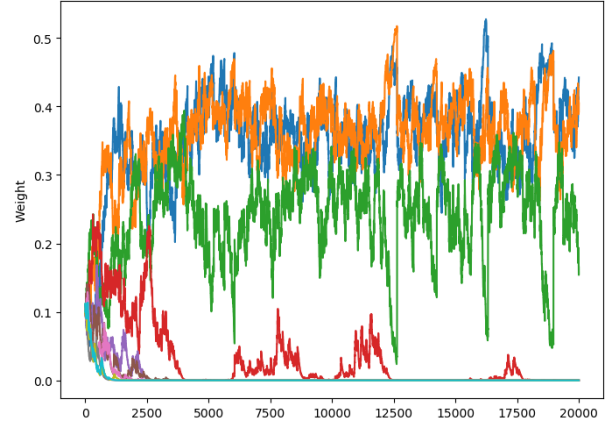
Figure 2a shows the regret of Exp3.M-VP versus the expected upper bound of the regret from Theorem 1. We can see that the actual regret of Exp3.M-VP has a smaller rate than its expected upper bound and the discrepancy becomes larger as time increases. Figure 2b shows the change of the normalized weight for each location over the entire time horizon. As shown in this figure, Exp3.M-VP chooses the top three locations (i.e. the blue, orange, and green curves) with the highest average reward only after a short period of time, and the rest of weights vanish to nearly 0. The reason why only three locations pop up is that M_t , i.e. the number of the arms played at each time, is within the set $\{1, 2, 3\}$. The fluctuations of the weights are partly due to the fact that the Exp3.M-VP algorithm needs to explore different locations in order to update the choice prediction and estimation, and partly due to the fact that the sum of the weights must always equal to M_t , which is changing over time.

B. Evaluations on Car-Hacking Dataset for the Defender

In this section we compare Exp3.M-VP with Exp3, Exp3.M, UCB, and the ϵ -greedy algorithms by implementing these algorithms over two in-vehicle network datasets from the Car-Hacking datasets. The Car-Hacking datasets are generated by logging the Controller Area Network (CAN) traffic via the OBD-II port from a real vehicle while message injection attacks were made. The Datasets each contain 300 intrusions of message injections over 26 unique CAN IDs. Each intrusion



(a) Exp3.M-VP regret (blue curve) and expected upper bound of regret (orange curve).



(b) Normalized weights of 10 arms over 20,000 time steps.

Fig. 2: Simulation of Exp3.M-VP on a ten-armed bandit problem.

is performed for 3 to 5 seconds, and each dataset has a total of 30 to 40 minutes of the CAN traffic. Specifically, we test the performance on the spoofing attack datasets, which were conducted on the RPM gauze and the driving gear. That is, among 26 arms representing CAN IDs, two of them (RPM gauze and driving gear) contained spoofing attacks.

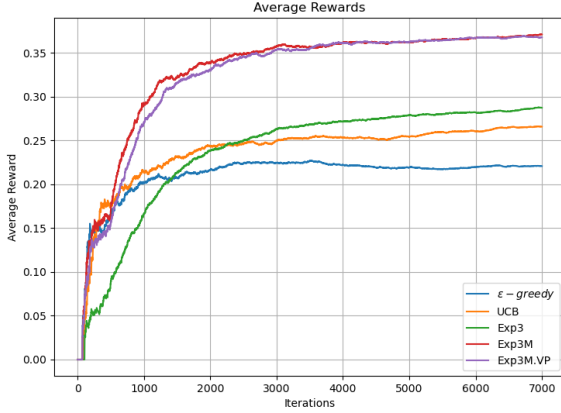


Fig. 3: Cumulative average rewards for ϵ -greedy, UCB, Exp3, Exp3.M, and Exp3.M-VP.

Figure 3 shows the cumulative average rewards for each bandit learning algorithm used by the defender. The experiments were conducted over $T = 7,000$ time steps, and the number of arms played by Exp3.M-VP was sampled from a truncated Gaussian distribution within the interval $[1, 3]$, with mean 2 and standard deviation 0.8. The number of arms played by Exp3.M was set to 3. We can see that both Exp3.M and Exp3.M-VP obtain higher cumulative average rewards than other single-play setting algorithms, due to the benefits from multiple or variable plays. Exp3.M-VP in this setting is a constrained version of Exp3.M, since the number of arms in Exp3.M (3) is an upper bound on the number of arms available to Exp3.M-VP ($[1, 3]$). This indicates that Exp3.M-VP may have access to a smaller number of arms due to resource constraints. To make it more challenging, Exp3.M-VP does not

know in advance the number of number of arms it may have access to in the future. Therefore, not surprisingly, Exp3.M obtains a slightly higher cumulative average reward than Exp3.M-VP. However, interestingly, eventually the cumulative average rewards of Exp3.M and Exp3.M-VP approach the same value. This demonstrates the power of the Exp3.M-VP algorithm: despite the fact that in average Exp3.M-VP plays fewer arms than Exp3.M, it can match the performance of Exp3.M. The reason is that only 2 out of 26 CAN-IDs contained spoofing attacks, and after a period of time (i.e. around 3500 iterations), both Exp3.M and Exp3.M-VP are able to identify the top two most rewarded CAN-IDs.

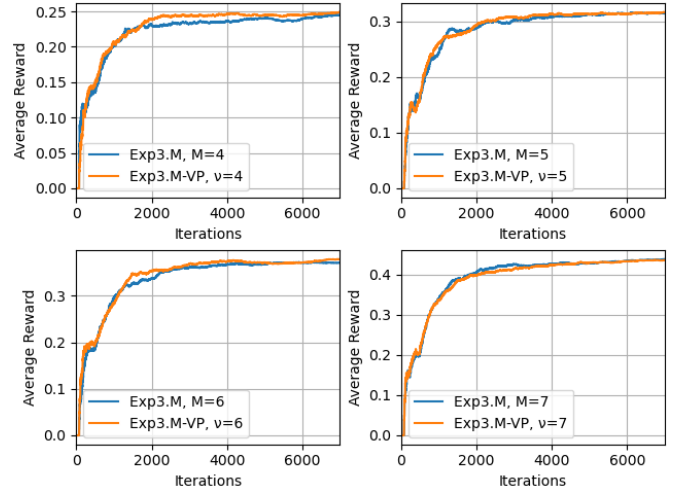


Fig. 4: Average reward of Exp3.M and Exp3.M-VP under different M_t and ν .

We further conduct sensitivity analysis on the number of arms played by Exp3.M and Exp3.M-VP. Specifically, we test the performance of the two algorithms with $M = \nu \in \{4, 5, 6, 7\}$, where M is the number of arms played by Exp3.M. For each ν , we sample M_t from a truncated Gaussian distribution within the interval $[\nu - 1, \nu + 1]$, with mean ν and standard deviation 0.8. As such, in this set of experiments the number of arms played by Exp3.M is the mean value of

the number of arms played by Exp3.M-VP. Figure 4 shows the results. This figure demonstrates the average reward of the two algorithms under four values for M and ν . We can see that the performance of the two algorithms are very close, mainly due to the fact that $M = \nu$. Note that here, in some instances Exp3.M-VP will have access to less resources/arms, and in some instances more. As a result, throughout the iterations, sometimes Exp3.M outperforms Exp3.M-VP, and sometimes it underperforms. However, eventually both algorithms reach the same reward and successfully identify the attacked arms. This again demonstrates the strength of Exp3.M-VP, because the number of arms are determined exogenously and therefore Exp3.M-VP is able to match the reward obtained by Exp3.M under uncertainty on the number of available arms at each time.

C. Simulations on Two Players

We now consider a game setting where two players, i.e., an attacker and a defender, are playing the pursuit-evasion game against each other. This corresponds to the realistic scenario where a malicious hacker is trying to compromise either the sensor/ECU in an in-vehicle sensor network, or the entire vehicle/infrastructure in an interconnected transportation system without being identified by the intrusion monitoring system. At the same time, the intrusion monitoring system is trying to identify as many compromised locations as possible to minimize the potential loss. We consider a ten-armed bandit problem for the two players, where the attacker adopts Exp3 and the defender adopts Exp3.M-VP. The scenario was simulated over $T = 100,000$ time steps, and the number of arms played by the defender was sampled from a truncated Gaussian distribution within the interval $[1, 3]$, with mean 2 and standard deviation 0.8. The parameter η for both Exp3 and Exp3.M-VP was set according to Corollary 1.1.

Figure 5 illustrates the average reward and the equilibrium reward for the two players. Since we have $N = 10$ and $\nu = 2$, according to Corollary 1.2 and Corollary 3.1, the equilibrium rewards for the attacker and the defender are 0.8 and 0.2, respectively. We can see that the average rewards of both players converge to the equilibrium rewards after a relatively short period, and after that the average rewards stay around the equilibrium reward with small fluctuations. The fluctuations are due to the fact the Exp3 and Exp3.M-VP use randomized policies and need to occasionally explore different locations in order to update the choice predictions and estimations.

IX. CONCLUSIONS

In this paper, we extend the adversarial/non-stochastic MPMAB to the case where the number of plays can change in time, and propose the Exp3.M-VP algorithm for obtaining the variable-play property. This extension is motivated by the uncertainty of resources allocated to the intrusion monitoring system to scan at each time in resource-constrained systems, such as an interconnected transportation system. We derive a sublinear regret bound for Exp3.M-VP, which simplifies to the existing bounds in the literature when the number of arms played at each time is constant. We introduce a game setting

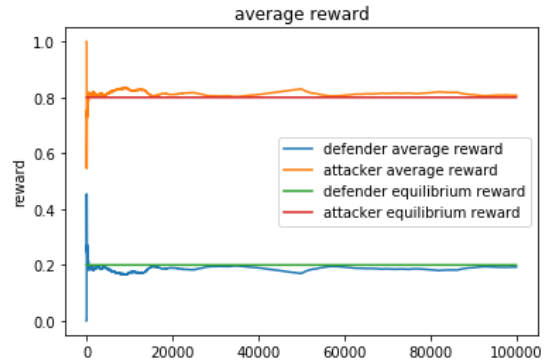


Fig. 5: Average reward of the attacker and the defender over 100,000 time steps.

where an attacker and a defender play a pursuit-evasion game against each other. The defender, who represents the intrusion monitoring system, adopts Exp3.M-VP and the attacker, who represents the malicious hacker, adopts Exp3. We derive the condition under which a Nash equilibrium of the strategic game exists. Finally, we consider heterogeneous rewards for arms, and obtain lower and upper bounds on the average rewards for the attacker in an infinite time horizon. We provide several numerical experiments that demonstrate our results.

This work provides insights on deploying an intrusion monitoring system either in an in-vehicle network or a transportation network: In order to minimize the potential loss of the system from cyber threats, one can either increase the average resources allocated to intrusion monitoring, or change the potential reward vector for each location to reduce the reward bound in Theorem 4. One of the potential extensions of this work is to consider the connectivity or correlations between different arms, which can take into account the spread of the cyber attacks, and use such information to facilitate the decision making of the intrusion monitoring system.

APPENDIX A HEDGE AND EXP3 ALGORITHMS

Algorithm 2 Hedge

- 1: Parameters: $\iota \in \mathbb{R}^+$.
- 2: Initialization: Set $r_k(1) := 0$ for all $k \in \mathcal{N}$.
- 3: **for** $t = 1, 2, \dots, T$ **do**
- 4: Choose action I_t according to the distribution

$$\beta_k = \frac{(1 + \iota)^{r_k(t)}}{\sum_{j=1}^N (1 + \iota)^{r_j(t)}}.$$

- 5: Receive the reward vector $x(t)$ and score gain $x_{I_t}(t)$.
- 6: Set $r_k(t+1) := r_k(t) + x_k(t)$ for all $k \in \mathcal{N}$.
- 7: **end for**

Algorithm 3 Exp3

-
- 1: Parameters: $\iota \in \mathbb{R}^+$ and $\eta \in [0, 1]$.
 - 2: Initialization: Initialize **Hedge**.
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Obtain the distribution vector $\beta(t) = (\beta_k(t), k \in \mathcal{N})$ from **Hedge**.
 - 5: Select action I_t to be k with probability

$$\hat{\beta}_k(t) = (1 - \eta)\beta_k(t) + \eta/N.$$

- 6: Receive the reward $x_{I_t}(t) \in [0, 1]$.
- 7: Return the simulated reward vector $\hat{x}(t) = (\hat{x}_k(t), k \in \mathcal{N})$ to **Hedge** with

$$\hat{x}_k(t) = \begin{cases} \frac{\eta}{N} \times \frac{x_{I_t}(t)}{\hat{\beta}_{I_t}(t)} & \text{if } k = I_t \\ 0 & \text{otherwise.} \end{cases}$$

8: **end for**

APPENDIX B
DEPROUND ALGORITHM

Algorithm 4 DepRound: The Dependent Rounding Algorithm

-
- 1: Inputs: Natural number $M < N$, marginal distribution $(p_k, k \in \mathcal{N})$ with $\sum_{k=1}^N p_k = M$
 - 2: Output: Subset \mathcal{N}_1 of \mathcal{N} such that $|\mathcal{N}_1| = M$
 - 3: **while** $\{k \in \mathcal{N} : 0 < p_k < 1\} \neq \emptyset$ **do**
 - 4: Choose distinct i and j such that $0 < p_i < 1$ and $0 < p_j < 1$
 - 5: Set $\rho = \min\{1 - p_i, p_j\}$ and $\zeta = \min\{p_i, 1 - p_j\}$
 - 6: Update p_i and p_j as

$$(p_i, p_j) = \begin{cases} (p_i + \rho, p_j - \rho) & \text{with probability } \frac{\zeta}{\rho + \zeta} \\ (p_i - \zeta, p_j + \zeta) & \text{with probability } \frac{\rho}{\rho + \zeta} \end{cases}$$

7: **end while**

8: **return** $\{k : p_k = 1, 1 \leq k \leq N\}$

APPENDIX C
PROOF OF THEOREM 1

Proof. Let $W_t := \sum_{k=1}^N w_k(t)$ and $W'_t := \sum_{k=1}^N w'_k(t)$. Then, at each time step t ,

$$\frac{W_{t+1}}{W_t} = \sum_{i \in S_0^c(t)} \frac{w_i(t+1)}{W_t} + \sum_{i \in S_0(t)} \frac{w_i(t+1)}{W_t} \quad (24a)$$

$$= \sum_{i \in S_0^c(t)} \frac{w_i(t)}{W_t} \exp\left(\frac{\eta M_t}{N} \hat{y}_i(t)\right) + \sum_{i \in S_0(t)} \frac{w_i(t)}{W_t} \quad (24b)$$

$$\leq \sum_{i \in S_0^c(t)} \frac{w_i(t)}{W_t} \left[1 + \frac{\eta M_t}{N} \hat{y}_i(t) + (e - 2) \times \left(\frac{\eta M_t}{N} \hat{y}_i(t)\right)^2 \right] + \sum_{i \in S_0(t)} \frac{w_i(t)}{W_t} \quad (24c)$$

$$= 1 + \frac{W'_t}{W_t} \sum_{i \in S_0^c(t)} \frac{w_i(t)}{W'_t} \left[\frac{\eta M_t}{N} \hat{y}_i(t) + (e - 2) \times \left(\frac{\eta M_t}{N} \hat{y}_i(t)\right)^2 \right] \quad (24d)$$

$$= 1 + \frac{W'_t}{W_t} \sum_{i \in S_0^c(t)} \frac{\frac{\hat{\alpha}_i(t)}{M_t} - \frac{\eta}{N}}{1 - \eta} \left[\frac{\eta M_t}{N} \hat{y}_i(t) + (e - 2) \times \left(\frac{\eta M_t}{N} \hat{y}_i(t)\right)^2 \right] \quad (24e)$$

$$\leq 1 + \frac{\eta}{(1 - \eta)N} \sum_{i \in S_0^c(t)} \hat{\alpha}_i(t) \hat{y}_i(t) + \frac{(e - 2)M_t \eta^2}{(1 - \eta)N^2} \times \sum_{i \in S_0^c(t)} \hat{\alpha}_i(t) \hat{y}_i^2(t) \quad (24f)$$

$$\leq 1 + \frac{\eta}{(1 - \eta)N} \sum_{i \in J_t \cap S_0^c(t)} y_i(t) + \frac{(e - 2)M_t \eta^2}{(1 - \eta)N^2} \times \sum_{i \in \mathcal{N}} \hat{y}_i(t). \quad (24g)$$

Inequality (24c) uses $e^a \leq 1 + a + a^2$, $\forall a \in [0, 1]$, equality (24e) holds because of step 11 in Algorithm 1, inequality (24f) uses the fact that $\frac{W'_t}{W_t} \leq 1$, and the last inequality (24g) holds because $\hat{\alpha}_i(t) \hat{y}_i(t) = y_i(t) \leq 1$ for $i \in J_t$ and $\hat{\alpha}_i(t) \hat{y}_i(t) = 0$ for $i \notin J_t$. Then, according to inequality (24g) and by summing over t , we have

$$\ln \frac{W_{T+1}}{W_1} = \sum_{t=1}^T \ln \frac{W_{t+1}}{W_t} \quad (25a)$$

$$\leq \sum_{t=1}^T \ln \left[1 + \frac{\eta}{(1 - \eta)N} \sum_{i \in J_t \cap S_0^c(t)} y_i(t) + \frac{(e - 2)M_t \eta^2}{(1 - \eta)N^2} \sum_{i \in \mathcal{N}} \hat{y}_i(t) \right] \quad (25b)$$

$$\leq \frac{\eta}{(1 - \eta)N} \sum_{t=1}^T \sum_{i \in J_t \cap S_0^c(t)} y_i(t) + \frac{(e - 2)b\eta^2}{(1 - \eta)N^2} \sum_{t=1}^T \sum_{i \in \mathcal{N}} \hat{y}_i(t). \quad (25c)$$

where inequality (25c) holds because $1 + y \leq e^y$ and $M_t \leq b$.

On the other hand, define \mathcal{A}_b^* as the best location index subset with b elements. Then,

$$\ln \frac{W_{T+1}}{W_1} \geq \ln \frac{\sum_{j \in \mathcal{A}_b^*} w_j(T+1)}{W_1} \quad (26a)$$

$$\geq \frac{\sum_{j \in \mathcal{A}_b^*} \ln w_j(T+1)}{b} - \ln \frac{N}{b} \quad (26b)$$

$$\geq \frac{\eta}{N} \sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0^c(t)} \hat{y}_j(t) - \ln \frac{N}{b}. \quad (26c)$$

where inequality (26a) holds because $\mathcal{A}_b^* \subseteq \mathcal{N}$, inequality (26b) comes from the inequality of arithmetic and geometric means, i.e. $\frac{1}{b} \sum_{j=1}^b y_j \geq \left(\prod_{j=1}^b y_j \right)^{\frac{1}{b}}$, and inequality (26c) is obtained by recursively applying step 15 of Algorithm 1, which results in equality (27):

$$w_j(T+1) = \exp \left((b\eta/N) \sum_{t: j \in S_0^c(t)} \hat{y}_j(t) \right). \quad (27)$$

Note that we also have

$$\sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0(t)} \hat{y}_j(t) \leq \sum_{t=1}^T \sum_{i \in S_0(t)} y_i(t) \quad (28a)$$

$$\leq \frac{1}{1-\eta} \sum_{t=1}^T \sum_{i \in S_0(t)} y_i(t) \quad (28b)$$

where inequality (28a) is due to the fact that $\hat{y}_j(t) = y_j(t), \forall j \in S_0(t)$, and the last inequality (28b) holds because $\eta \in (0, 1]$.

Combining (25c), (26c), (28a), and (28b), we have:

$$\sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0^c(t)} \hat{y}_j(t) + \sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0(t)} \hat{y}_j(t) - \frac{N}{\eta} \ln \frac{N}{b} \quad (29a)$$

$$\leq \frac{1}{(1-\eta)} G_{\text{Exp3.M-Vp}}^J(T) + \frac{(e-2)\eta b}{(1-\eta)N} \sum_{t=1}^T \sum_{i \in \mathcal{N}} \hat{y}_i(t) \quad (29b)$$

Taking expectations of both sides of inequality (29), we obtain

$$\sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0^c(t)} \hat{y}_j(t) + \sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0(t)} \hat{y}_j(t) - \frac{N}{\eta} \ln \frac{N}{b} \quad (30a)$$

$$\leq \frac{1}{(1-\eta)} E [G_{\text{Exp3.M-Vp}}^J(T)] + \frac{(e-2)\eta b}{(1-\eta)N} \sum_{t=1}^T \sum_{i \in \mathcal{N}} y_i(t) \quad (30b)$$

$$\leq \frac{1}{(1-\eta)} E [G_{\text{Exp3.M-Vp}}^J(T)] + \frac{(e-2)\eta b}{(1-\eta)a} G_{\max}(T), \quad (30c)$$

where inequality (30b) uses the fact that $E[\hat{y}_i(t) | S(1), \dots, S(t-1)] = y_i(t)$, and

$$\sum_{t=1}^T \sum_{i \in \mathcal{N}} y_i(t) \leq \frac{N}{a} G_{\max}(T). \quad (31)$$

Since $\mathcal{A}_b^* = \cup_{M_t} \mathcal{A}_{M_t}^*$ trivially holds, we have

$$G_{\max}(T) - \frac{N}{\eta} \ln \frac{N}{b} \leq \sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0^c(t)} \hat{y}_j(t) \quad (32)$$

$$+ \sum_{j \in \mathcal{A}_b^*} \sum_{t: j \in S_0(t)} \hat{y}_j(t) - \frac{N}{\eta} \ln \frac{N}{b} \quad (33)$$

Therefore, by combining (30) and (32), we obtain the inequality stated in the Theorem 1. ■

APPENDIX D

PROOF OF COROLLARY 1.1

Proof. For any $T > 0$, we have $G_{\max}(T) \leq bT$. If $bT \leq \sqrt{\frac{Na \ln(N/b)}{a+(e-2)b}}$, then the bound is trivial since the expected regret cannot be more than bT . Otherwise, by Theorem 1, the expected regret is at most

$$2\sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \sqrt{bTN \ln \frac{N}{b}}$$

by plugging in $\eta = \sqrt{\frac{Na \ln(N/b)}{a+(e-2)bT}}$. ■

APPENDIX E

PROOF OF THEOREM 3

Proof. Note that

$$\bar{r}_{\infty} = 1 - \liminf_{T \rightarrow \infty} E \left[\frac{1}{T} G_{\text{Exp3.M-Vp}}^J(T) \right] \quad (34a)$$

$$\leq 1 - \liminf_{T \rightarrow \infty} \frac{1}{T} \left(G_{\max}(T) - 2\sqrt{\left(1 + (e-2)\frac{b}{a}\right)} \sqrt{bTN \ln \frac{N}{b}} \right) \quad (34b)$$

$$= 1 - \liminf_{T \rightarrow \infty} \frac{1}{T} G_{\max}(T) \quad (34c)$$

$$\leq \frac{N-a}{N} \quad (34d)$$

for any policy π of the attacker, where the last inequality (34d) comes from the fact that $G_{\max}(T) \geq \frac{Ta}{N}$ for any defender's policy γ .

Under the greedy policy we have $\hat{\beta}_{g(t)}(t) \leq \frac{b}{N}$, which implies $r(t) \geq \frac{N-b}{N}$ for any t . Therefore by using the greedy policy π_{greedy} , we have $\bar{r}_{\infty} \geq \frac{N-b}{N}$. ■

APPENDIX F

PROOF OF LEMMA 5

The proof of Lemma 5 is an extension of the proof of Lemma 4 in [39]. The main difference is that the matrix \mathbf{H} is now an $N \times |\mathcal{C}(\mathcal{N}, a)|$ matrix compared to the one in the original proof which is $N \times N$.

Proof. 1) The problem (18) is equivalent to

$$\max_{d \in \Delta_N} \min_{u \in \Delta_N} \sum_{n=1}^{|\mathcal{C}(\mathcal{N}, a)|} \sum_{k=1}^N \left(\mu_k d_k - \sum_{j \in J_n} \mu_j d_j \right) u_n \quad (35)$$

which can be rewritten in matrix form:

$$\max_{d \in \Delta_N} \min_{u \in \Delta_N} \mathbf{d}^\top \mathbf{H} \mathbf{u}, \quad (36)$$

where \top denotes the transpose, and

$$\mathbf{H} = \begin{bmatrix} 0, & \mu_1, & \dots & \mu_1 \\ & \dots & & \\ 0, & \mu_a, & \dots & \mu_a \\ \mu_{a+1}, & \dots & & \mu_{a+1} \\ & \dots & & \\ \mu_{N-a+1}, & \dots & & 0 \\ & \dots & & \\ \mu_N, & \dots & & 0 \end{bmatrix}$$

where each column j represents one set $\mathcal{S} \subseteq \mathcal{C}(\mathcal{N}, a)$ such that for all $i \in \mathcal{S}$, $H_{ij} = 0$ and for all $i \in \mathcal{N} \setminus \mathcal{S}$, $H_{ij} = \mu_i$. The remaining proof is the same as the original proof.

Now consider a zero-sum game with the payoff matrices for the row and the column players being \mathbf{H} and $-\mathbf{H}$, whose mixed strategy vectors are \mathbf{d} and \mathbf{u} , respectively. Any optimal solution \mathbf{d}^* to the problem (18) is a Nash equilibrium strategy for the row player, and by the indifference condition, we obtain for any $j \in \text{supp}(\mathbf{d}^*)$,

$$\sum_{k \neq j} k \in \text{supp}(\mathbf{d}^*) = \text{Const.}, \quad (37)$$

which implies $\mu_k d_k^* = \mu_j d_j^*$ for any $k, j \in \text{supp}(\mathbf{d}^*)$.

2) The second part of the Lemma is proved by contradiction. Assume that there exist $i \in \text{supp}(\mathbf{d}^*)$ and $j \in \mathcal{N} \setminus \text{supp}(\mathbf{d}^*)$ such that $\mu_j > \mu_i$. Let ϵ be a constant such that $\epsilon = \mu_k d_k^*$ for any $k \in \text{supp}(\mathbf{d}^*)$. Then consider a feasible solution \mathbf{d} , where $d_k = 0$ for all $k \in ((\mathcal{N} \setminus \text{supp}(\mathbf{d}^*)) \setminus \{j\}) \cup \{i\}$, and $d_k = d_k^* + \epsilon$ for all $k \in (\text{supp}(\mathbf{d}^*) \setminus \{i\}) \cup \{j\}$, with $\epsilon = d_i^*(1 - \mu_i/\mu_j)/K^*$, which yields a higher objective value. ■

APPENDIX G PROOF OF LEMMA 6

Proof. Consider the following linear program:

$$\text{minimize}_{l, p} p \quad (38a)$$

$$\text{s.t. } p + \mu_k l_k \leq \mu_k, \quad (38b)$$

$$\sum_k l_k \leq b, \quad (38c)$$

$$l_k \leq 1, \quad (38d)$$

$$l_k \geq 0. \quad (38e)$$

It is easy to see that problem (23) is lower bounded by the problem (38).

Then the dual of the program (38) can be written as

$$\text{maximize}_{d, q} \sum_{k=1}^N \mu_k d_k - q \quad (39a)$$

$$\text{s.t. } \sum_{k=1}^N \mu_k d_k \leq \frac{Nq}{b}, \quad (39b)$$

$$\sum_{k=1}^N d_k = 1, \quad (39c)$$

$$d_k \geq 0. \quad (39d)$$

Note that program (39) is equivalent to the following problem

$$\text{maximize}_{d \in \Delta_N} \sum_{k=1}^N \mu_k d_k - \max_{J \in \mathcal{C}(N, b)} \sum_{j \in J} \mu_j d_j, \quad (40)$$

which is essentially problem (18), except for changing the set $\mathcal{C}(\mathcal{N}, a)$ to $\mathcal{C}(\mathcal{N}, b)$. Therefore problem (40) has the optimal value $\frac{K^* - b}{\sum_{k=1}^{K^*} \mu_k}$, which provides us with the lower bound. ■

REFERENCES

- [1] J. G. Ibanez, S. Zeadally, and J. Contreras-Castillo, "Integration challenges of intelligent transportation systems with connected vehicle, cloud computing, and internet of things technologies," *IEEE Wireless Communications*, vol. 6, no. 22, pp. 122–128, 2015.
- [2] Y. Cao, C. Xiao, B. Cyr, Y. Zhou, W. Park, S. Rampazzi, Q. A. Chen, K. Fu, and Z. M. Mao, "Adversarial sensor attack on lidar-based perception in autonomous driving," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 2267–2281.
- [3] Y. Cao, C. Xiao, D. Yang, J. Fang, R. Yang, M. Liu, and B. Li, "Adversarial objects against lidar-based autonomous driving systems," *arXiv preprint arXiv:1907.05418*, 2019.
- [4] Y. Feng, S. Huang, Q. A. Chen, H. X. Liu, and Z. M. Mao, "Vulnerability of traffic control system under cyberattacks with falsified data," *Transportation research record*, vol. 2672, no. 1, pp. 1–11, 2018.
- [5] J. Petit and S. E. Shladover, "Potential cyberattacks on automated vehicles," *IEEE Transactions on Intelligent transportation systems*, vol. 16, no. 2, pp. 546–556, 2014.
- [6] S. Checkoway, D. McCoy, B. Kantor, D. Anderson, H. Shacham, S. Savage, K. Koscher, A. Czeskis, F. Roesner, T. Kohno *et al.*, "Comprehensive experimental analyses of automotive attack surfaces," in *USENIX Security Symposium*, vol. 4. San Francisco, 2011, pp. 447–462.
- [7] A. Weimerskirch and R. Gaynier, "An overview of automotive cybersecurity: Challenges and solution approaches," in *TrustED@ CCS*, 2015, p. 53.
- [8] C. Yan, W. Xu, and J. Liu, "Can you trust autonomous vehicles: Contactless attacks against sensors of self-driving vehicle," *DEF CON*, vol. 24, no. 8, p. 109, 2016.
- [9] S. Parkinson, P. Ward, K. Wilson, and J. Miller, "Cyber threats facing autonomous and connected vehicles: Future challenges," *IEEE transactions on intelligent transportation systems*, vol. 18, no. 11, pp. 2898–2915, 2017.
- [10] E. Fok, "An introduction to cybersecurity issues in modern transportation systems," *ITE Journal*, vol. 3, p. 19, 2013.
- [11] K. B. Kelarestaghi, K. Heaslip, M. Khalilikhah, A. Fuentes, and V. Fesmann, "Intelligent transportation system security: hacked message signs," *SAE International Journal of Transportation Cybersecurity and Privacy*, vol. 1, no. 11-01-02-0004, pp. 75–90, 2018.
- [12] F. van Wyk, Y. Wang, A. Khojandi, and N. Masoud, "Real-time sensor anomaly detection and identification in automated vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 3, pp. 1264–1276, 2019.
- [13] Y. Wang, N. Masoud, and A. Khojandi, "Real-time sensor anomaly detection and recovery in connected automated vehicle sensors," *IEEE Transactions on Intelligent Transportation Systems*, 2020.

- [14] —, “Anomaly detection in connected and automated vehicles using an augmented state formulation,” *arXiv preprint arXiv:2004.09496*, 2020.
- [15] M. Marchetti, D. Stabili, A. Guido, and M. Colajanni, “Evaluation of anomaly detection for in-vehicle networks through information-theoretic algorithms,” in *2016 IEEE 2nd International Forum on Research and Technologies for Society and Industry Leveraging a better tomorrow (RTSI)*. IEEE, 2016, pp. 1–6.
- [16] M. Muter and N. Asaj, “Entropy-based anomaly detection for in-vehicle networks,” in *2011 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2011, pp. 1110–1115.
- [17] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, “Gambling in a rigged casino: The adversarial multi-armed bandit problem,” in *Proceedings of IEEE 36th Annual Foundations of Computer Science*. IEEE, 1995, pp. 322–331.
- [18] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [19] A. György, T. Linder, G. Lugosi, and G. Ottucsák, “The on-line shortest path problem under partial monitoring,” *Journal of Machine Learning Research*, vol. 8, no. Oct, pp. 2369–2403, 2007.
- [20] N. Nisan, T. Roughgarden, E. Tardos, and V. V. Vazirani, *Algorithmic Game Theory*. Cambridge University Press, 2007.
- [21] J. D. Abernethy, E. Hazan, and A. Rakhlin, “Competing in the dark: An efficient algorithm for bandit linear optimization,” in *21st Annual Conference on Learning Theory - COLT 2008, Helsinki, Finland, July 9-12, 2008*, R. A. Servedio and T. Zhang, Eds. Omnipress, 2008, pp. 263–274.
- [22] N. Cesa-Bianchi and G. Lugosi, “Combinatorial bandits,” *Journal of Computer and System Sciences*, vol. 78, no. 5, pp. 1404–1422, 2012.
- [23] V. Anantharam, P. Varaiya, and J. Walrand, “Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays-part i: Iid rewards,” *IEEE Transactions on Automatic Control*, vol. 32, no. 11, pp. 968–976, 1987.
- [24] R. Agrawal, M. Hegde, and D. Teneketzis, “Multi-armed bandit problems with multiple plays and switching cost,” *Stochastics and Stochastic reports*, vol. 29, no. 4, pp. 437–459, 1990.
- [25] J. Komiyama, J. Honda, and H. Nakagawa, “Optimal regret analysis of thompson sampling in stochastic multi-armed bandit problem with multiple plays,” *arXiv preprint arXiv:1506.00779*, 2015.
- [26] Y. Xia, T. Qin, W. Ma, N. Yu, and T.-Y. Liu, “Budgeted multi-armed bandits with multiple plays,” in *IJCAI*, 2016, pp. 2210–2216.
- [27] T. Uchiya, A. Nakamura, and M. Kudo, “Algorithms for adversarial bandit problems with multiple plays,” in *Algorithmic Learning Theory*, M. Hutter, F. Stephan, V. Vovk, and T. Zeugmann, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 375–389.
- [28] E. Fouché, J. Komiyama, and K. Böhm, “Scaling multi-armed bandit algorithms,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1449–1459.
- [29] A. Lesage-Landry and J. A. Taylor, “The multi-armed bandit with stochastic plays,” *IEEE Transactions on Automatic Control*, vol. 63, no. 7, pp. 2280–2286, 2017.
- [30] H. Sedjelmaci, M. Hadji, and N. Ansari, “Cyber security game for intelligent transportation systems,” *IEEE Network*, vol. 33, no. 4, pp. 216–222, 2019.
- [31] C. Kiekintveld, M. Jain, J. Tsai, J. Pita, F. Ordóñez, and M. Tambe, “Computing optimal randomized resource allocations for massive security games,” in *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*, 2009, pp. 689–696.
- [32] A. Sinha, T. H. Nguyen, D. Kar, M. Brown, M. Tambe, and A. X. Jiang, “From physical security to cybersecurity,” *Journal of Cybersecurity*, vol. 1, no. 1, pp. 19–35, 2015.
- [33] T. Alpcan and S. Buchegger, “Security games for vehicular networks,” *IEEE Transactions on Mobile Computing*, vol. 10, no. 2, pp. 280–290, 2010.
- [34] M. N. Mejri, N. Achir, and M. Hamdi, “A new security games based reaction algorithm against dos attacks in vanets,” in *2016 13th IEEE Annual Consumer Communications & Networking Conference (CCNC)*. IEEE, 2016, pp. 837–840.
- [35] S. Bahamou, M. D. El Ouadghiri, and J.-M. Bonnin, “When game theory meets vanet’s security and privacy,” in *proceedings of the 14th international conference on advances in mobile computing and multi media*, 2016, pp. 292–297.
- [36] H. Sedjelmaci, S. M. Senouci, and N. Ansari, “Intrusion detection and ejection framework against lethal attacks in uav-aided networks: A bayesian game-theoretic methodology,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 5, pp. 1143–1153, 2016.
- [37] R. Vidal, O. Shakernia, H. J. Kim, D. H. Shim, and S. Sastry, “Probabilistic pursuit-evasion games: theory, implementation, and experimental evaluation,” *IEEE transactions on robotics and automation*, vol. 18, no. 5, pp. 662–669, 2002.
- [38] V. Navda, A. Bohra, S. Ganguly, and D. Rubenstein, “Using channel hopping to increase 802.11 resilience to jamming attacks,” in *IEEE INFOCOM 2007-26th IEEE International Conference on Computer Communications*. IEEE, 2007, pp. 2526–2530.
- [39] Q. Wang and M. Liu, “Learning in hide-and-seek,” *IEEE/ACM transactions on networking*, vol. 24, no. 2, pp. 1279–1292, 2015.
- [40] R. Gandhi, S. Khuller, S. Parthasarathy, and A. Srinivasan, “Dependent rounding and its applications to approximation algorithms,” *Journal of the ACM (JACM)*, vol. 53, no. 3, pp. 324–360, 2006.
- [41] S. Bubeck, N. Cesa-Bianchi *et al.*, “Regret analysis of stochastic and nonstochastic multi-armed bandit problems,” *Foundations and Trends® in Machine Learning*, vol. 5, no. 1, pp. 1–122, 2012.
- [42] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, “The non-stochastic multiarmed bandit problem,” *SIAM journal on computing*, vol. 32, no. 1, pp. 48–77, 2002.
- [43] P. Auer, N. Cesa-Bianchi, and P. Fischer, “Finite-time analysis of the multiarmed bandit problem,” *Machine learning*, vol. 47, no. 2-3, pp. 235–256, 2002.
- [44] R. S. Sutton, A. G. Barto *et al.*, *Introduction to reinforcement learning*. MIT press Cambridge, 1998, vol. 135.
- [45] E. Seo, H. M. Song, and H. K. Kim, “Gids: Gan based intrusion detection system for in-vehicle network,” in *2018 16th Annual Conference on Privacy, Security and Trust (PST)*, Aug 2018, pp. 1–6.