

Scalable Game-Focused Learning of Adversary Models: Data-to-Decisions in Network Security Games

Kai Wang, Andrew Perrault, Aditya Mate, Milind Tambe
Harvard University
{kaiwang,aperrault,aditya_mate}@g.harvard.edu,milind_tambe@harvard.edu

ABSTRACT

Previous approaches to adversary modeling in network security games (NSGs) have been caught in the paradigm of first building a full adversary model, either from expert input or historical attack data, and then solving the game. Motivated by the need to disrupt the multibillion dollar illegal smuggling networks, such as wildlife and drug trafficking, this paper introduces a fundamental shift in learning adversary behavior in NSGs by focusing on the accuracy of the model using the downstream game that will be solved. Further, the paper addresses technical challenges in building such a game-focused learning model by i) applying graph convolutional networks to NSGs to achieve tractability and differentiability and ii) using randomized block updates of the coefficients of the defender’s optimization in order to scale the approach to large networks. We show that our game-focused approach yields scalability and higher defender expected utility than models trained for accuracy only.

KEYWORDS

Adversarial multi-agent learning; Game theory for practical applications

ACM Reference Format:

Kai Wang, Andrew Perrault, Aditya Mate, Milind Tambe. 2020. Scalable Game-Focused Learning of Adversary Models: Data-to-Decisions in Network Security Games. In *Proc. of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2020)*, Auckland, New Zealand, May 9–13, 2020, IFAAMAS, 9 pages.

1 INTRODUCTION

Many real-world security problems present the challenge of how to allocate limited resources to large number of important targets, including infrastructure [13], transportation systems [32], urban crime [47], and web security [42]. *Stackelberg security games* (SSGs) are frequently used to study the interaction between defender and attacker and optimally allocate the security resources accordingly. *Network security games* (NSGs) [10, 37, 43], a natural extension of SSGs, describe a strategic adversarial interaction between an attacker and a defender on a graph. The attacker’s goal is to take a path from a starting location to a target without being caught by the defender. The defender declares (i.e., attacker surveils) a mixed strategy of how she will deploy her security resources to the edges of the network. NSGs are relevant in many real-world settings such as wildlife conservation [8, 24], infrastructure protection [20], and nuclear material smuggling [28, 33].

One key challenge in applying NSGs in the real world is learning an adversary’s behavior from historical data. Past works [1, 5, 29] in security games have shown that constructing bounded rationality adversary models from data greatly improves performance of deployed models because attackers often behave quite differently from how rational models would suggest. A particular motivation for this paper is wildlife smuggling [9, 36, 48], a natural NSG domain where large amounts of historical attack data is available in the form of past seizures.

Almost all previous work on security games approaches the problem of adversary modeling by first building a full adversary model that aims to predict the adversary behavior as accurately as possible [2, 7, 8, 31]. In early work, the judgments of human experts were used to estimate the adversary’s preferences [39]. Later, in domains where historical attack data was available, machine learning was used to construct models instead (starting from Letchford et al. [23]). In NSGs, building an adversary model to maximize accuracy has several key limitations. First, the model is selected without any consideration of the impact of errors downstream. Prediction errors on paths that are frequently taken by the adversary have a large impact on defender utility, but are weighted the same as errors on paths that are rarely taken. Secondly, standard adversary models require human feature engineering to apply to NSGs due to a great variety of paths from the attacker’s starting location to each potential target [12, 15, 16, 45]. Once the adversary model is determined, the following defender utility maximization problem can be solved by any optimization techniques, including bilevel optimization [24], branch and cut [11], and double oracle [20].

Our approach represents a fundamental shift: we take an end-to-end, game-focused approach, focusing on learning a model that yields a high defender utility. More specifically, we take the downstream defender utility maximization problem into account while learning the adversary model. To that end, we use a graph convolutional neural network architecture to learn the adversary’s behavior, which allows us to overcome both of the issues of prior work. First, assuming we can differentiate through the defender’s optimization problem, we can train the entire model end-to-end because the predictive model is differentiable, i.e., to take the optimization problem into account while training. Second, the graph convolutional network automatically extracts features from the graph, meaning that hand engineering is not necessary. Nevertheless, several challenges must be overcome to implement this approach: principally, poor scalability of naive end-to-end training and non-convexity of the game-focused objective.

A summary of our contributions is as follows: first, we construct a *graph convolution*-based adversary model for NSGs. This model is fully differentiable, does not require manual selection of path features, and transmits target value information over the network.

Proc. of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2020), B. An, N. Yorke-Smith, A. El Fallah Seghrouchni, G. Sukthakar (eds.), May 9–13, 2020, Auckland, New Zealand. © 2020 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

Second, we develop a randomized block update scheme for differentiating through optimization problems, whose computation time is usually more than quadratic in terms of the number of variables due to the computation of Hessian matrix and matrix inversion. Such computational issue is especially influential for optimization problems with a huge number of variables, which is commonly seen in NSGs as every edge corresponds to one individual decision variable. In these cases, randomized block update can largely reduce the time complexity. We further provide an approximation guarantee relative to the complete derivatives, and we show empirically that our approach greatly improves scalability. We also show that through judicious use of the standard predictive loss as regularization, we can escape local minima in the end-to-end loss function.

Related Work. There is a rich literature on learning adversary behavior models in Stackelberg security games (SSGs) (starting from Letchford et al. [23]), but learning in NSGs has received much less attention. While SSGs generalize NSGs, the scalability concerns are quite different because reducing NSGs to SSGs may create exponentially many targets—one for each path to the target in the NSG. Thus, applying standard attacker bounded rationality models, such as quantal response (QR) [25, 26] and subjective utility quantal response (SUQR) [31] is nontrivial. Yang et al. [45] and Ford et al. [12] reduced NSGs to SSGs by considering each individual path as an attacker pure strategy. Their approach scales poorly, creating exponentially many paths in many networks. It also relies on hand-crafting path features that capture adversary behavior well. Other authors have developed models that use Markovian dynamics to model the attacker. Gutfraind et al. [15] and Abbasi et al. [2] assume the attacker does not receive any information beyond the neighboring nodes—attackers do not make any decisions that are more long term than a single timestep. Gutfraind et al. [16] takes the opposite approach: attackers follow a path that minimizes some cost (such as the risk of being caught) with randomness in the individual decisions. This adds some global information, but requires the model designer to specify the choice of cost function in advance.

Past work in adversary modeling in SSGs has viewed the problem of constructing an adversary model and solving the defender’s optimization as completely separate problems and does not consider the impact of errors in the defender model on the quality of the optimization outcome, with a few exceptions. Sinha et al. [38] and Haghtalab et al. [17] relate the predictive accuracy of the learned model to the defender’s expected utility. In the case of Haghtalab et al., this view motivates the use of a non-standard loss function to achieve better utility. However, even these papers take a fundamentally two-stage approach: the model is trained independently of any information about the game itself, such as the defender’s utilities. Perrault et al. [35] takes a game-focused approach to SSGs, but the issues that arise in NSGs are different and require a greater focus on scalability.

A major challenge in our work is differentiating through the nonconvex defender optimization problem. Recent work has developed general approaches for differentiating convex problems [3]. Perrault et al. [35] present an approach for a limited class of nonconvex problems. Our setting is challenging in two ways. First,

we have a decision variable for each edge in the network and these approaches scale poorly (more than quadratically) in the number of variables. Second, our setting is more severely nonconvex than that of Perrault et al.

2 BACKGROUND

Stackelberg Security Games. A *Stackelberg security game (SSG)* [39, 46] is a two-player sequential game. The defender aims to protect a set of targets T with limited budget b which can only protect up to b targets. Each target $t \in T$ is associated with a defender penalty $U^d(t) \leq 0$ and an attacker reward $U^a(t) \geq 0$ when the target is successfully attacked. For simplicity, we assume there is no reward and penalty when the attacker is caught or fails to reach the target. Once the defender commits to her mixed strategy, the attacker can conduct surveillance to observe the defender’s mixed strategy and choose one target to attack accordingly. We denote the defender’s mixed strategy by $\mathbf{x} \in \mathbb{R}^{|T|}$, where $0 \leq \mathbf{x}_t \leq 1$ denotes the marginal probability that target t is protected. The budget constraint can be written as $\mathbf{1}^\top \mathbf{x} \leq b$. On the attacker side, we use $\mathbf{q}(\mathbf{x}, \xi)$ to represent the attacker’s behavior, where $\mathbf{q}_t(\mathbf{x}, \xi)$ (or \mathbf{q}_t if there is no ambiguity) is the probability of attacking target t , and ξ is the available features revealed to both the defender and the attacker, e.g., the attacker payoff value $U^a(t) \forall t \in T$ can be considered as a feature. Notice that \mathbf{q} is a function of the defender strategy \mathbf{x} and the feature ξ , which implies that the attacker can be reactive to the defender strategy and select the target based on the underlying feature. We can write the defender’s utility function as:

$$\text{DefU}(\mathbf{x}; \mathbf{q}) = \sum_{t \in T} \mathbf{q}_t(\mathbf{x}, \xi) U^d(t) (1 - \mathbf{x}_t). \quad (1)$$

This includes the case where the attacker is fully rational, where $\mathbf{q}_t(\mathbf{x}, \xi) = 1$ if $t = \arg \max_{t' \in T} (1 - \mathbf{x}_{t'}) U^a(t')$ else 0.

Bounded Rationality in SSGs. *Quantal response (QR)* [25] models the attacker’s behavior by setting the probability that each target is attacked to be proportional to the exponential of its payoff scaled by a constant. *Subjective utility quantal response (SUQR)* [31], which fits data better than QR in practice, sets the probability proportional to the exponential of a subjective utility or an attractiveness function of the attacker:

$$\mathbf{q}_t(\mathbf{x}, \xi) \propto \exp(-\omega \mathbf{x}_t + \Phi(t, \xi)), \quad (2)$$

where $\omega > 0$ is a constant representing the attacker’s risk aversion and $\Phi(t, \xi)$ denotes the subjective utility of target t given feature ξ .

Network Security Games. *Network security games (NSGs)* [10, 30] are SSGs played on a graph structure. Given an undirected (or directed) graph $G = (V, E)$, the defender allocates a limited number of checkpoints along edges in E , while the attacker tries to find a path from a source to a target without being caught. We divide the set of all vertices V into targets $T = \{t_1, t_2, \dots, t_{|T|}\}$ and non-targets $S = \{s_1, s_2, \dots, s_{|S|}\}$ (or potential sources). At each time, the attacker appears in one potential source $s \in S$ drawn from a given prior distribution $\pi \in \mathbb{R}^{|S|}$. From the defender’s perspective, the defender strategy $\mathbf{x}_e \forall e \in E$ is the marginal probability of covering edge e . Similarly, the defender has a limited number of resources b to protect the targets.

We use $\alpha = \{v_1, v_2, \dots, v_{|\alpha|}\}$ to denote a path which starts from a source $v_1 \in S$ and ends with a target $v_{|\alpha|} \in T$. We use \mathcal{A}

to denote the set of all possible paths from any source to any target, which could be exponentially many or infinitely many when the graph contains any cycle. Similar to SSGs, let $U^d(t)$ be the defender’s payoff when the target t is attacked successfully and U^d_{caught} be the defender’s payoff when the attacker is caught. Let $U^d = \{U^d(t_1), \dots, U^d(t_{|T|}), U^d_{\text{caught}}\} \in \mathbb{R}^{|T|+1}$ denote the defender’s payoff vector. In addition, we assume each node $v \in \mathcal{V}$ has a node feature vector $\xi_v \in \mathbb{R}^D$ consisting of characteristics of node v , e.g., the attacker payoff of the current node $U^a(v)$ if $v \in T$. We use $\xi \in \mathbb{R}^{|V| \times D}$ to denote all the node features in graph G .

Bounded Rationality in NSGs. In this paper, we assume the attacker to be boundedly rational, where the attacker’s behavior is characterized by a function $\mathbf{q}(\mathbf{x}, \xi)$, where $\mathbf{q}_\alpha(\mathbf{x}, \xi)$ represents the probability of choosing path α under coverage \mathbf{x} and feature ξ . Given the coverage \mathbf{x} , we can compute the defender expected utility:

$$\text{DefU}(\mathbf{x}; \mathbf{q}) = \sum_{\alpha \in \mathcal{A}} \mathbf{q}_\alpha(\mathbf{x}, \xi) U^d(\alpha) \prod_{e \in \alpha} (1 - \mathbf{x}_e), \quad (3)$$

where $U^d(\alpha) = U^d(t)$ is the defender utility when the attacker successfully passes through α to attack its target t .

The difference between Equation 1 and 3 is that there are multiple layers of protection along the path α . Therefore the probability of successfully attacking a target is the product of all the success probabilities of crossing each edge e in the path. The defender’s optimization problem is generally hard. For example, if the function $\mathbf{q}(\mathbf{x}, \xi)$ is given by full rationality restricted to only polynomial many paths \mathcal{A} , the defender optimization problem is NP-hard [20]. Furthermore, the set of all possible paths \mathcal{A} could be exponentially large or infinitely many when there is any cycle.

Graph Convolutional Networks. There has been much recent attention paid to *graph convolutional networks (GCNs)* [18, 21, 27]. Given a graph, the convolutional layers in GCNs can transmit information through message passing, which allows information to propagate to distant nodes and be aggregated in a non-linear fashion. GCNs are much more expressive than hand-crafted features. In this paper, we apply GCNs, parameterized by θ , to map each node $v \in V$ and the entire node features ξ with graph structure to a scalar $\Phi(v, \xi; \theta)$, which represents the extent that the attacker is “pulled” toward that node. The message passing in GCNs is similar to the information gathering conducted by the adversary, where a rough understanding of faraway targets is available to the adversary.

3 ADVERSARY MODEL

Our attacker model is Markovian—the probability of using a path α can be decomposed into the product of transition probabilities:

$$\mathbf{q}_\alpha(\mathbf{x}, \xi) = \prod_{e \in \alpha} \mathbf{q}_e(\mathbf{x}, \xi). \quad (4)$$

Motivated by the SUQR model, we propose a **local SUQR** model, which assumes the probability that the attacker moves from u to v using edge $e = (u, v)$ is proportional to $\exp(-\omega \mathbf{x}_{u \rightarrow v} - \eta \mathbf{y}_v + \Phi(v, \xi; \theta)) \forall v \in N_{\text{out}}(u)$. $\Phi(v, \xi; \theta)$ represents the subjective utility or attractiveness of node v parameterized by θ , which can be learned by GCN. \mathbf{y}_v , with a weight $\eta \geq 0$, represents the downstream future risk or coverage perceived by the attacker at node v . In other words, the attacker tends to move toward the target with higher

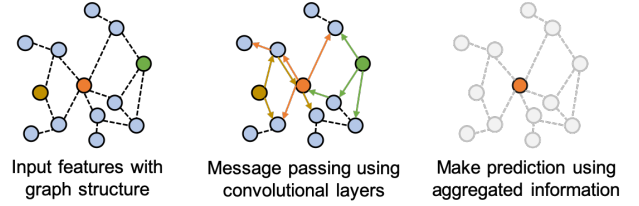


Figure 1: The convolutional layers of GCNs can propagate and aggregate information in a non-linear fashion. In NSGs, such message passing ability corresponds to the attacker’s ability of conducting surveillance to neighbor nodes.

attractiveness $\Phi(v, \xi; \theta)$, but avoids using the edge $e = (u, v) \in E$ with higher coverage $\mathbf{x}_{u \rightarrow v}$ and avoids moving towards nodes v with higher future risk \mathbf{y}_v .

Given a defender coverage strategy, there are many heuristic ways to obtain a measure of future risk. For example, we can follow the above Markovian behavior without the effect of the future risk, where the probability of being caught can be analytically computed efficiently. Another heuristic is the shortest distance to any target, as suggested by Gutfraind et al. [16]. The only restriction put on the choice of the future risk is differentiability.

We can compute the transition probability from u to any $v \in N_{\text{out}}(u)$ as:

$$\mathbf{q}_{u \rightarrow v}(\mathbf{x}, \xi; \theta) = \frac{\exp(-\omega \mathbf{x}_{u \rightarrow v} - \eta \mathbf{y}_v + \Phi(v, \xi; \theta))}{\sum_{v' \in N_{\text{out}}(u)} \exp(-\omega \mathbf{x}_{u \rightarrow v'} - \eta \mathbf{y}_{v'} + \Phi(v', \xi; \theta))}. \quad (5)$$

Unlike previous boundedly rational models [12, 45], we do not need to enumerate all the feasible paths, which could be exponentially large. Unlike the nonreactive Markovian model [15], our model is reactive to the defender’s strategy. Unlike Gutfraind et al. [16], we are not limited to noisily following a shortest path.

In local SUQR, the path structure is automatically encoded in the reactive Markovian behavior. Since the edge coverage effect is involved in the transition probability, the probability of taking a path is also exponentially proportional to the total coverage along the path, which is also included in other bounded rational models [12, 45]. The flexibility and the generalizability of the attractiveness function allow us to apply any graph learning algorithms to extract the adversary behavior. Compared to previous hyperparameters tuning models, our model is more expressive and can adapt to a broader range of adversary behavior.

4 PROBLEM STATEMENT

For each instance, a directed graph $G = (V, E)$ with node features ξ is presented to both the defender and the attacker. The attacker has a hidden rationality function \mathbf{q}^* , which is a function of node features ξ and the defender coverage \mathbf{x} . The defender first chooses a coverage $\{\mathbf{x}_e\}_{e \in E}$ under the budget constraint $\mathbf{1}^\top \mathbf{x} \leq b$. The attacker observes \mathbf{x} and then behaves based on his own rationality function \mathbf{q}^* . We assume that the defender has access to historical play between the defender and the attacker, which can be used to form an estimate of the adversary behavior. The goal of the defender is to maximize the received expected reward.

5 TWO-STAGE LEARNING FOR NETWORK SECURITY GAMES

The main comparison of the remainder of the paper is between our GCN-based adversary model implemented as two-stage vs. our game-focused methods. Thus, we briefly describe the two-stage approach that we consider.

Predictive Model. A two-stage approach fits the GCN-based attractiveness function $\Phi(v, \xi) \forall v \in V$ to minimize the difference between predicted behavior \mathbf{q} given by Equation 5 and the corresponding true attacker behavior \mathbf{q}^* . Given the attacker behavior \mathbf{q}^* and a prediction \mathbf{q} , we can define the loss by either matrix norm or the KL-divergence of the path distribution inferred by two behaviors under previous coverage \mathbf{x} and features ξ . These losses are generally infeasible to compute since there are infinite many possible paths. In practice, however, we often have paths sampled from the true behavior \mathbf{q}^* we can use to approximately compute the KL-divergence between two behaviors. Given the choice of loss function \mathcal{L} , we can train a model \mathbf{q} by minimizing the average loss:

$$E_{(\mathbf{x}, \xi, \mathbf{q}^*) \in D} \mathcal{L}(\mathbf{q}^*, \mathbf{q}; \mathbf{x}, \xi) \quad (6)$$

Prescriptive Model. Given a graph G , node features ξ , and predicted attacker behavior \mathbf{q} , the defender’s goal is to choose an optimal coverage \mathbf{x}^* satisfied the budget constraint to maximize her own objective value.

When the defender strategy \mathbf{x} is chosen, the attacker follows his own Markovian behavior $\mathbf{q}(\mathbf{x}, \xi)$. But due to the allocated coverage, the attacker will be caught with probability \mathbf{x}_e when he passes through edge e . This can be cast as an absorbing Markov chain, where the probability of crossing an edge e is $\mathbf{q}_e(\mathbf{x}, \xi)(1 - \mathbf{x}_e)$, and the rest of the probability the attacker will be caught and turned into a dummy caught state v_{caught} . We also assume that once the attacker reaches either any terminal or caught state v_{caught} , the attacker cannot go back to any other states, i.e., these are absorbing states. Therefore, given a coverage \mathbf{x} , we can model the attacker’s behavior as an absorbing Markov chain. We can analytically compute the corresponding defender utility. To align with the standard minimization formulation, we denote the *negative* defender utility by $f(\mathbf{x}, \mathbf{q})$. For ease of notation, we omit the presence of node features. The optimization problem is given by:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}, \mathbf{q}) \\ \text{s.t.} \quad & \mathbf{1}^\top \mathbf{x} \leq b, \quad 0 \leq \mathbf{x}_e \leq 1 \quad \forall e \in \mathcal{E} \end{aligned} \quad (7)$$

Unfortunately, the function f is neither convex nor submodular when the attacker is reactive. The standard approach is to apply constrained black-box optimization solvers to solve the problem, e.g., Sequential Least Squares Programming (SLSQP) [6, 22].

6 NAIVE GAME-FOCUSED LEARNING FOR NETWORK SECURITY GAMES

In general, a good predictive model does not necessarily imply a high defender utility in the second stage. Sometimes a slightly inaccurate prediction might lead to a better final decision. This happens frequently especially when the predictive model cannot perfectly represent the ground truth. For example, in our case, the model relies on the Markovian assumption and SUQR assumption in

Algorithm 1: Naive Game-focused Learning [35]

```

1 Input: Training data  $D$ , initialized  $\text{GCN}(\cdot, \cdot; \theta) : V \times \xi \rightarrow \mathbb{R}$ 
2 while until converge do
3   for  $(G, \mathbf{q}^*, \xi) \in D$  do
4     Compute prediction  $\mathbf{q}$  in Eq. 5 by  $\Phi = \text{GCN}(V, \xi; \theta)$ 
5     Find optimum  $\mathbf{x}^{\text{opt}}$  of Optimization 7
6      $Q = \frac{\partial^2 f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}, p = \frac{\partial f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}} - Q\mathbf{x}^*$ 
7     Re-solve QP:  $\mathbf{x}^* = \underset{\mathbf{x} \text{ feasible}}{\text{argmin}} \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \mathbf{x}^\top p$ 
8     Update  $\theta$  by gradient  $\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}^*} \frac{d\mathbf{x}^*}{dp} \frac{dp}{d\theta}$ 
9 Return: trained model  $\text{GCN}(\cdot, \cdot; \theta)$ 

```

Equation 5, which might not be able to fully recover the underlying attacker behavior.

Game-focused learning, instead, can directly optimize the final solution quality by back-propagating from the final solution quality all-the-way back to the predictive model. Game-focused learning has been proven to be able to outperform a standard two-stage learning approach [35], finding a shortcut to better final solution quality. However, the major issue of back-propagation is the non-differentiable optimization layer in the prescriptive state. Amos et al. [4] provides a method to differentiate through the optimization layer when the optimization program is convex; Perrault et al. [35] instead used quadratic function as a surrogate to deal with the case when the optimization program is non-convex.

More specifically, the idea of tackling non-convex function in Perrault et al. [35] is to approximate the non-convex function by a quadratic function around a local minimum \mathbf{x}^{opt} using Taylor expansion, which can be written as:

$$f(\mathbf{x}, \mathbf{q}) \approx f(\mathbf{x}^{\text{opt}}, \mathbf{q}) + (\Delta \mathbf{x})^\top \frac{\partial f}{\partial \mathbf{x}} + \frac{1}{2} (\Delta \mathbf{x})^\top \frac{\partial^2 f}{\partial \mathbf{x}^2} (\Delta \mathbf{x}) \quad (8)$$

where $\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}^{\text{opt}}$. They use this approximate quadratic program (QP) as a surrogate of the non-convex optimization problem, where the optimal solution \mathbf{x}^* of QP matches the local optimum \mathbf{x}_{opt} computed before. This allows us to differentiate through a QP and compute the gradient of optimal solution \mathbf{x}^* with respect to the linear coefficient $p = \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}$.

$$\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\theta} = \frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}^*} \frac{d\mathbf{x}^*}{dp} \frac{dp}{d\theta} \quad (9)$$

where $p = \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}$ is a function of \mathbf{q} with $\frac{dp}{d\theta} = \frac{dp}{dq} \frac{dq}{d\Phi} \frac{d\Phi}{d\theta}$ can be decomposed and computed. Equation 9 gives us the gradient of the final solution quality with respect to the model parameter θ , which allows us to directly run stochastic gradient descent end-to-end. We apply this approach to our domain. The algorithm is sketched in Algorithm 1 and Figure 2(b).

Issues of Game-focused Learning. Although game-focused learning ideally can achieve better final performance compared to two-stage learning, in this section, we point out two main issues that arise when this game-focused learning is applied to NSGs: scalability and non-convexity.

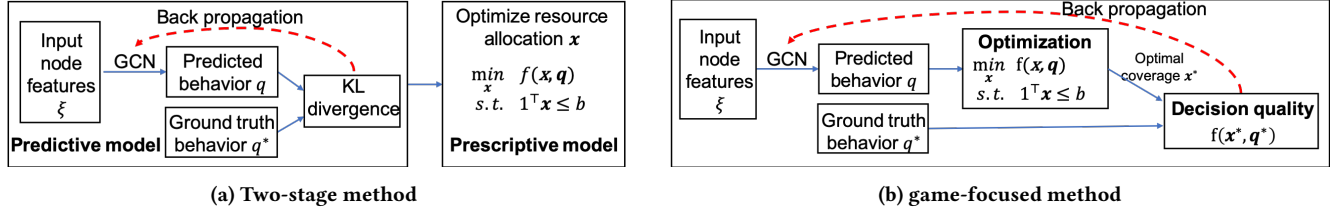


Figure 2: Two-stage method trains the behavior model by minimizing the predictive loss, while the game-focused method trains the behavior model by optimizing the final decision quality.

- *Scalability:* In the forward and backward paths of solving QP (Equation 8), we need to solve and be able to back-propagate through the QP, which involves the computation of matrix inverse. Taking matrix inverse grows between quadratic and cubically as the size of the decision variable \mathbf{x} grows. Moreover, in order to compute the Taylor expansion 8, we need to compute the Hessian $\frac{\partial^2 f}{\partial \mathbf{x}^2}$ explicitly, which is usually the major bottleneck of the computation cost when the target function f is complex.

- *Non-convexity:* In the non-convex setting, the objective function $f(\mathbf{x}, \mathbf{q})$ can be non-convex in both \mathbf{x} and \mathbf{q} . The gradient-based approaches rely on updating model parameters θ and thus \mathbf{q} to improve the solution quality. However, since the f is non-convex in \mathbf{q} , it could create non-convex searching space for gradient-based approaches, which could easily get stuck in local optimum or saddle points. Two-stage methods escape this problem because their loss function $\mathcal{L}(\mathbf{q}, \mathbf{q}^*)$ in Equation 6 is convex, which gradient-based approaches can more easily handle.

7 IMPROVING NAIVE GAME-FOCUSED LEARNING

In this section, we provide a scalable randomized block update approach to resolve the scalability issue, which also suggests a block game-focused algorithm as a scalable version of game-focused learning approach. To resolve the non-convexity issue, we apply the intermediate loss as a regularization, which helps game-focused methods escape local minimums. We further provide theoretical guarantees to link the randomized block update to the naive game-focused learning approach.

7.1 Block Game-focused Learning

Instead of using the entire Taylor expansion (Equation 8) to approximate the objective function locally, we can use a partial Taylor expansion with respect to a subset of variables to approximate it:

$$f(\mathbf{x}, \mathbf{q}) \approx f(\mathbf{x}^*, \mathbf{q}) + (\Delta \mathbf{x}_C)^T \frac{\partial f}{\partial \mathbf{x}_C} + \frac{1}{2} (\Delta \mathbf{x}_C)^T \frac{\partial^2 f}{\partial \mathbf{x}_C^2} (\Delta \mathbf{x}_C), \quad (10)$$

where $C \subset \{1, 2, \dots, |E|\}$ is a subset of indices and \mathbf{x}_C is the corresponding truncation over indices C of the entire variables \mathbf{x} . Equation 10 is equivalent to freezing the variables outside of C and applying Taylor expansion to the rest of them. In this formulation, we only need to compute the Hessian with respect to \mathbf{x}_C . When the size of C is significantly smaller than the original variable size $|E|$, it can save the computational time of Hessian quadratically. Furthermore, while back-propagating through the KKT conditions,

Algorithm 2: Block Game-focused Learning

- 1 **Input:** Training data D , initialized $\text{GCN}(\cdot, \cdot; \theta) : V \times \xi \rightarrow \mathbb{R}$, block size k
- 2 **while** until converge **do**
- 3 **for** $(G, \mathbf{q}^*, \xi) \in D$ **do**
- 4 Compute prediction \mathbf{q} in Eq. 5 by $\Phi = \text{GCN}(V, \xi; \theta)$
- 5 Find optimum \mathbf{x}^{opt} of Optimization 7
- 6 Sample $C \subset \{1, 2, \dots, |E|\}$ with $|C| = k$
- 7 $Q_{CC} = \frac{\partial^2 f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}_C^2} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}, p_C = \frac{\partial f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}_C} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}} - Q_{CC} \mathbf{x}_C^*$
- 8 Re-solve QP: $\mathbf{x}_C^* = \underset{\mathbf{x}_C \text{ feasible}}{\text{argmin}} \frac{1}{2} \mathbf{x}_C^T Q_{CC} \mathbf{x}_C + \mathbf{x}_C^T p_C$
- 9 Update θ by gradient $\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}_C^*} \frac{d\mathbf{x}_C^*}{dp_C} \frac{dp_C}{d\theta}$
- 10 **Return:** trained model $\text{GCN}(\cdot, \cdot; \theta)$

the QP formulation of Equation 10 results in a smaller size of quadratic term, which can reduce the computation of matrix inverse. The block-wise chain rule can be written as:

$$\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\theta} \approx \frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}_C^*} \frac{d\mathbf{x}_C^*}{dp_C} \frac{dp_C}{d\theta} \quad (11)$$

where $p = \frac{\partial f}{\partial \mathbf{x}_C} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}$, $\frac{dp_C}{d\theta} = \frac{dp_C}{dq} \frac{dq}{d\Phi} \frac{d\Phi}{d\theta}$. When the block size is smaller, the approximation can be more inaccurate. But we will show in the later section that the block gradient is an approximation to the entire gradient.

All the above reasons suggest a randomized block update algorithm, which is described in Algorithm 2. The algorithm randomly samples a block of variables to compute Hessian and back-propagate accordingly. In comparison, Algorithm 1 requires to compute the entire Hessian matrix $Q = \frac{\partial^2 f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}$, which is usually very expensive. Instead, Algorithm 2 only requires the computation of a block Hessian $Q_{CC} = \frac{\partial^2 f(\mathbf{x}, \mathbf{q})}{\partial \mathbf{x}_C^2} \Big|_{\mathbf{x}=\mathbf{x}^{\text{opt}}}$, which can save at least quadratic amount of Hessian computation depending on the block size. It can also reduce the running time of the following quadratic program due to reducing the number of variables.

7.2 Block Selection

In Algorithm 2, the idea of block game-focused learning is to restrict the focus to a subset of variables and to update accordingly. The choice of the sampled block could affect the convergence rate. Here we propose three block selection approaches: i) *random* approach

selects block uniformly at random; ii) *coverage*-based approach randomly selects indices with probability proportional to \mathbf{x}^* , which guarantees that there is space for the variables in the block to reallocate coverage; iii) *derivative*-based approach selects indices with probability proportional to the magnitude of the derivatives $\frac{df(\mathbf{x}^*, \mathbf{q})}{dx_i^*}$, which is the weight put on the chain rule.

7.3 Regularization

Another issue associated with the naive game-focused learning method is the non-convex objective function, where gradient-based approaches can encounter issues of local optimums and saddle points. Instead, the two-stage approach optimizes the intermediate loss, which is generally convex in the prediction space. Therefore, we propose to add a weighted two-stage loss as a regularization to smoothify the final objective value. As the training epochs increase, the weight put on the two-stage loss drops exponentially with a decay rate 0.95, pulling the learning back to game-focused methods. This regularization technique helps resolve the non-convexity issue of naive game-focused method, which can achieve better performance afterward.

7.4 Approximation Guarantees

In this section, we will show that both Algorithm 1 and 2 have 0 gradient when the prediction perfectly matches to the ground truth, showing that both algorithms are stable at the global optimum. Later on, we will show that Algorithm 2 is an approximate version of Algorithm 1. This shows that our block game-focused approach can not only achieve scalability due to the reduction in Hessian and QP computation, but it is also aligned with the standard naive game-focused approach with theoretical guarantees.

THEOREM 7.1. *When the intermediate prediction matches the ground truth, i.e., $\mathbf{q}(\cdot, \cdot; \theta^*) = \mathbf{q}^*$, we have $\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\theta} |_{\theta=\theta^*} = 0$ for both Algorithm 1 and Algorithm 2 with any block C .*

This theorem implies that if the predictive model is rich enough and able to reach the ground truth, then the gradient computed in both algorithms is equal to 0 at the ground truth. So if we can avoid getting stuck by local optimum, then both algorithms will be able to learn the ground truth. This is also true for the two-stage learning when the loss is defined as any convex norms.

THEOREM 7.2. *The quadratic programs in Algorithm 1 and Algorithm 2 share the same primal solutions on the block C . They also share the same dual solution on the non-degenerate constraints containing at least one variable in the block.*

When restricting to variables inside the block, there are some degenerate constraints containing only variables outside of the block, which are always satisfied in the block QP. Thus, there is no restriction put on the dual variable corresponding to these degenerative constraints, which we have no control on them. But in this theorem, we prove that the dual solution to the other valid constraints will match to the dual solution given by the QP in Algorithm 1.

THEOREM 7.3. *Given the primal solution \mathbf{x}^* and the dual solution λ^* of the quadratic program in Algorithm 1 with linear constraints G, h, A, b , the Hessian $Q = \frac{\partial^2 f}{\partial \mathbf{x}^2}$, linear coefficient $p = \frac{\partial f}{\partial \mathbf{x}}$, and the*

sampled indices $C \subset \{1, 2, \dots, |E|\}$, the gradient $\frac{d\mathbf{x}_C^}{dp_C} \in \mathbb{R}^{|C| \times |C|}$ computed in Algorithm 2 is an approximation to the block component of the gradient $\frac{d\mathbf{x}^*}{dp} \in \mathbb{R}^{|E| \times |E|}$ computed in Algorithm 1. More specifically,*

$$\left\| \left(\frac{d\mathbf{x}^*}{dp} \right)_{CC} - \frac{d\mathbf{x}_C^*}{dp_C} \right\| \leq \frac{\Delta + \Delta_C}{\mu_{\min}(Q)} \max(\lambda^*, 1) \|K_{CC}^{-1}\| \left\| \left(\frac{d\mathbf{x}^*}{dp} \right)_{CC} \right\| \quad (12)$$

where $\Delta = \|G^T G + A^T A\|$, $\Delta_C = \|Q_{CC}^{-1} Q_{CC}^{-1}\|$, and $\mu_{\min}(Q)$ is the smallest eigenvalue of positive definite matrix Q . K_{CC} is the KKT matrix given by the quadratic program in Algorithm 2.

The Δ in the numerator is a constant that only depends on the constraint matrices. The other term Δ_C depends on the choice of block C , which measures the magnitude of the off-diagonal elements of the Hessian matrix Q . This is usually a small term when the Hessian Q is diagonally dominant. Another interesting finding is that this bound depends on the convexity of the Hessian Q . When the Hessian is more convex, then the smallest eigenvalue of Q is also larger, giving a stronger bound in Theorem 7.3. The last term K_{CC}^{-1} measures the stability of the KKT matrix K_{CC} . We can get a good bound if the KKT matrix K_{CC} is far from singular. Greif et al. [14] provides various bounds on the eigenvalues of the KKT matrix. However, in general, poor constraints can still lead to a KKT matrix close to singular. It also indicates that a good choice of C can imply a more stable KKT matrix, leading to a better estimate in Theorem 7.3.

Theorem 7.3 also implies an alternative explanation to Algorithm 2, where the gradient in Algorithm 2 is an approximation to the partial gradient with indices C in Algorithm 1:

$$\frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}_C^*} \left(\frac{d\mathbf{x}^*}{dp} \right)_{CC} \frac{dp_C}{d\theta} \approx \frac{df(\mathbf{x}^*, \mathbf{q}^*)}{d\mathbf{x}_C^*} \frac{d\mathbf{x}_C^*}{dp_C} \frac{dp_C}{d\theta} \quad (13)$$

which implies that Algorithm 2 can be thought as an approximate block-wise gradient descent of Algorithm 1, which relates to the literature of block coordinate gradient descent [34, 41].

8 EXPERIMENTS

In this section, we compare *two-stage (TS)*, *naive game-focused (naive-GF)* mentioned in Section 6, *block game-focused (block-GF)*, and *regularized block game-focused (reg-block-GF)* methods on synthetic data to show that our block game-focused and regularized block game-focused methods can achieve better performance especially in larger instances. These two methods are also able to scale up to large instances, where the naive game-focused method cannot. Lastly, we study the convergence and scalability of the block game-focused and regularized block game-focused methods with different block sizes and block sampling methods. This allows us to choose the right block size to balance between solution quality and scalability.

8.1 Synthetic Data Generation

8.1.1 Graph and features: we first randomly generate a graph G with various node sizes, 5 random sources with uniform initial distribution π , and 5 random targets with defender penalties $U^d(t) \forall t \in T$ drawn from $[-10, -5]$ uniformly at random. We focus on stochastic block model [19] and geometric graphs [44], which

can respectively model community structures and physical road networks¹. For each node in the graph, we draw an attractiveness value, depending on the shortest distance to the targets plus a uniform noise, as the attacker’s unbiased preference. We also randomly generate the past coverage \mathbf{x} subject to budget constraints. To generate the node features ξ , we feed the private attractiveness values to a randomly initialized GCN, where the GCN will output a fixed size vector per node as our node features ξ . A different level of Gaussian noise was added to the features to model the noise in the real-world scenario.

8.1.2 Attacker behavior: we choose $\omega = 4$ as the risk aversion parameter suggested by Perrault et al. [35] and Abbasi et al. [1], and set $\eta = 0$ to ignore the future risk factor for the sake of simplicity. For each instance with given attractiveness and the defender coverage, we simulate 100 attacks by initializing the attacker at one of the sources and following the localized SUQR behavior described in Section 3 until the attacker reaches to one of the targets. These sampled paths Λ are used to reconstruct a Markovian behavior: $\mathbf{q}_{u \rightarrow v}^*(\mathbf{x}, \xi) := \frac{|\{e=(u,v), e \in \alpha, \alpha \in \Lambda\}|}{|\{e=(u,w), e \in \alpha, \alpha \in \Lambda, w \in N(u)\}|}$ [40], which is then used as our ground truth to evaluate the solution quality². Each instance is composed of the graph G , past coverage \mathbf{x} , node features ξ , the attacker behavior \mathbf{q}^* , and the sampled paths Λ (only used in two-stage method).

8.1.3 Training, validating, and testing: we generate 50 instances ($G, \mathbf{q}^*, \mathbf{x}^*, \xi$) as our entire dataset, which are randomly separated into training, validating, testing set with size 35, 5, 10. The model is trained on the training set for 100 epochs, where the best model is chosen from the 100 epochs with the highest score in the validation set. In the following experiments, to achieve statistical significance, for every method and different setup, we ran 50 independent trials and recorded the average results on the testing set.

8.2 Solution Quality

In this section, we compare the solution quality of all methods on stochastic block models and geometric graphs. We generate a set of random graphs with features as described in Section 8.1.1, where Gaussian noise with std. of 0.2 is added to the features to model noisy real-world data. We set $b = 2$. As our goal is efficient approaches for adversary models in large-scale NSGs, the focus of this paper is then on experimenting with many different settings (graph sizes and types), techniques (different variations of game-focused learning), noise, and other variables in building an adversary model. In addition, since we care more about how much defender utility that various learning approaches can improve, we focus on the *counterfactual regret*, which is defined as the gap between the defender utility of our solution and the true optimum

¹For stochastic block model, we separate nodes into communities with 10 nodes in each community, then connect nodes within the same community with probability 0.4 and nodes not in the same community with probability 0.1. For geometric graph, we randomly places nodes in a unit square and connects nodes with distance smaller than 0.2.

²The reason of using sampled paths instead of the actual generated attractiveness values as our ground truth is to align with the real-world data, where it is almost impossible to have access to the underlying attacker preference or Markovian behavior; instead, we generally have access to the paths or edges where illegal activities have been found, which can be used as sampled paths or edges and used to reconstruct the Markovian behavior as we did here.

when the ground truth is given in advance. Smaller regret implies that the solution is closer to the actual optimum.

In Figure 3(a) and 3(b), we can see that our regularized block game-focused method outperforms two-stage method (note that all of the improvements in the average regret reported by the reg-block-GF method over the two-stage method are statistically significant with $p < 0.05$). When the instance gets larger, the difference between two approaches also gets larger, showcasing the limit of the standard two-stage behavior learning approach. In Figure 4(a) and 4(b), we compare the solution quality of different game-focused methods. Due to the computational issue, the naive game-focused method can only scale up to graphs with 40 nodes. The block game-focused method can scale up to larger instances but it sacrifices some solution quality compared to the naive game-focused approach. Finally, the regularized block game-focused method can achieve both scalability and solution quality by using the block update and regularization term.

8.3 The Impact of Noise

Figure 5(a) and 5(b) compare the performances under different level of noise, where a noise with std. of r is added to the normalized features. We can see that the more noise implies larger regret and poorer performance. But we can also notice that the gap between regularized block game-focused method and the two-stage method gets larger when more noise is introduced. This is probably due to the mismatch between the low intermediate loss and the good final solution quality when the feature is noisy. This also explains why regularized block game-focused method can outperform two-stage in Figure 3 when the features are noisy.

8.4 Scalability

Figure 6(a) and 6(b) show the scalability of all game-focused methods. We limit the training time to be up to 48 hours. Any programs last more than that were cut and the corresponding results were recorded. Naive game-focused method can only handle graphs with up to 40 nodes and it scales extremely poorly. Our proposed methods, block game-focused and regularized block game-focused with a block size $\#nodes/2$, can scale to much larger instances.

8.5 Block Size Selection

To study the effect of block size, we select various block sizes proportional to the total number of variables and run the block game-focused learning and regularized block game-focused methods to compare the convergence. In Figure 7(a), we can see that for the block-game-focused method, the convergence and the final performance are better when the block size is larger. Figure 7(b) shows the convergence of regularized block game-focused method with different block sizes. In this case, a larger block size still helps, but the difference is relatively tiny.

Figure 7(c) shows the running time of the forward (lines 4-5) and backward path (lines 6-9 in Algorithm 2) for the block game-focused method with various block sizes, where forward path solves prescriptive stage with black-box optimization and the backward path requires computing the Hessian and solving the quadratic program to back-propagate. In practice, we would like to select a block size such that the running time of forward and backward paths are of

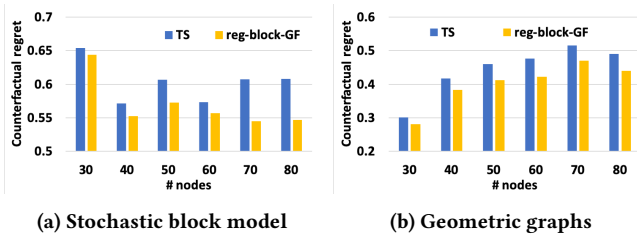


Figure 3: Solution quality comparison between two-stage and regularized block game-focused method. The difference in solution quality gets larger when the graph size increases.

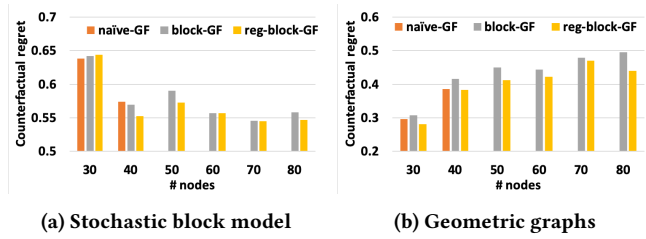


Figure 4: Solution quality comparison between game-focused methods. Randomized block update can improve scalability while the regularization can improve the solution quality.

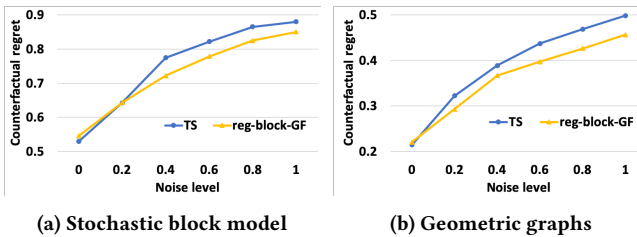


Figure 5: The figures show the effect of noise to all the methods, where regularized block game-focused method is more resilient to noise in the features.

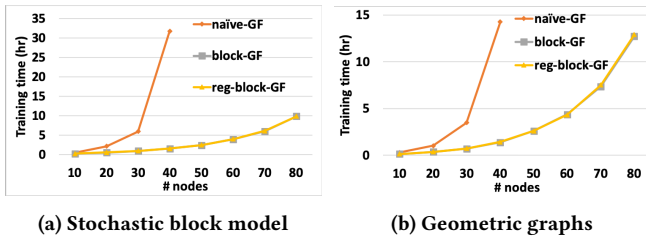


Figure 6: Naive game-focused method can only scale up to 40 nodes. Instead, block game-focused and regularized block game-focused can solve larger instances with 80 nodes.

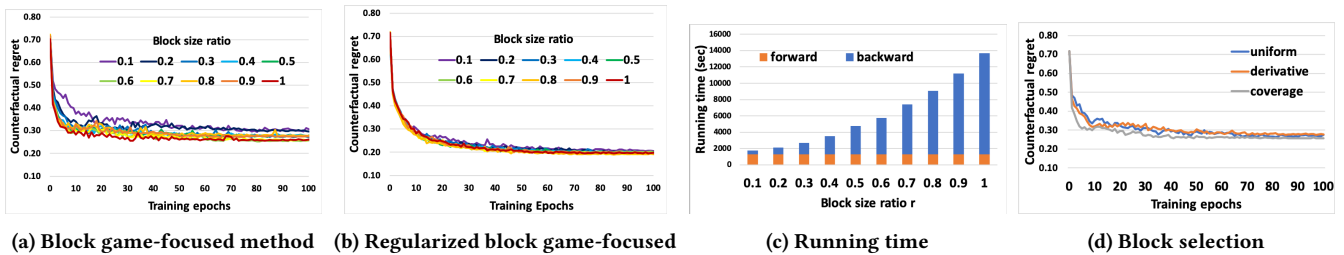


Figure 7: Figure (a) and (b) show the convergence rate of different block sizes. Figure (c) shows the running time of backward path for different block sizes, which grows significantly more than linear. Figure (d) shows the effect of different block sampling methods. All methods converge with slightly different speed, where coverage-based sampling is the best and it is also what we use in other experiments.

the same order to balance between the convergence and scalability, which explains the reason that we eventually choose block size = #nodes/2 for all other experiments. Lastly, Figure 7(d) compares different block selections mentioned in Section 7.2, where convergence speed differs but mostly lead to the same point. Coverage-based selection converges the most quickly, and thus we use it throughout the other experiments.

9 CONCLUSIONS

In this paper, we introduce a fundamentally different behavior learning approach, game-focused learning, to network security games, placing the downstream defender utility maximization problem into the loop of behavior learning. We propose a novel local SUQR model as our adversary model, where GCNs can be applied to automatically handle the information propagation in the graph.

We further identify two existing issues of game-focused learning method: scalability and non-convexity, which are addressed by our block game-focused and by regularizing respectively. Block game-focused method can largely reduce the computational cost while maintaining the focus on the final solution quality as naive game-focused learning does. We also provide theoretical guarantees on the block game-focused method. In the experimental section, we run extensive experiments to verify the reduction on the training time and show an improvement in terms of solution quality. The block game-focused method reduces the training time, but sacrifices a little solution quality, while regularized block game-focused can achieve both speed and performance.

Acknowledgments. This research was supported by MURI Grant Number W911NF-17-1-0370 and W911NF-18-1-0208.

REFERENCES

- [1] Yasaman Abbasi, Debarun Kar, Nicole Sintov, Milind Tambe, Noam Ben-Asher, Don Morrison, and Cleotilde Gonzalez. 2016. Know Your Adversary: Insights for a Better Adversarial Behavioral Model. In *CogSci*.
- [2] Yasaman Dehghani Abbasi, Martin Short, Arunesh Sinha, Nicole Sintov, Chao Zhang, and Milind Tambe. 2015. Human adversaries in opportunistic crime security games: evaluating competing bounded rationality models. In *Proc. of Advances in Cognitive Systems*.
- [3] Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. 2019. Differentiable Convex Optimization Layers. In *NeurIPS-19*. Vancouver.
- [4] Brandon Amos and J. Zico Kolter. 2017. OptNet: Differentiable optimization as a layer in neural networks. In *ICML-17*. Sydney.
- [5] Jana Arsovska and Panos A Kostakos. 2008. Illicit arms trafficking and the limits of rational choice theory: the case of the Balkans. *Trends in Organized Crime* 11, 4 (2008), 352–378.
- [6] Dimitri P. Bertsekas and John. N. Tsitsiklis. 1996. *Neuro-dynamic Programming*. Athena, Belmont, MA.
- [7] Sarah Cooney, Kai Wang, Elizabeth Bondi, Thanh Nguyen, Phebe Vayanos, Hailey Winetrobe, Edward A Cranford, Cleotilde Gonzalez, Christian Lebiere, and Milind Tambe. 2019. Learning to Signal in the Goldilocks Zone: Improving Adversary Compliance in Security Games. In *ECMLPKDD-19*. Würzburg.
- [8] Fei Fang, Peter Stone, and Milind Tambe. 2015. When Security Games Go Green: Designing Defender Strategies to Prevent Poaching and Illegal Fishing. In *IJCAI-15*. Buenos Aires, 2589–2595.
- [9] Peyton Ferrier et al. 2009. *The economics of agricultural and wildlife smuggling*. Technical Report. Springer.
- [10] Matteo Fischetti, Ivana Ljubić, Michele Monaci, and Markus Sinnl. 2016. *Interdiction games and monotonicity*. Technical Report. Technical Report, DEI, University of Padova.
- [11] Matteo Fischetti, Ivana Ljubić, Michele Monaci, and Markus Sinnl. 2019. Interdiction games and monotonicity, with application to knapsack problems. *INFORMS Journal on Computing* (2019).
- [12] Benjamin Ford, Thanh Nguyen, Milind Tambe, Nicole Sintov, and Francesco Delle Fave. 2015. Beware the soothsayer: From attack prediction accuracy to predictive reliability in security games. In *GameSec-15*. 35–56.
- [13] Jiarui Gan, Bo An, and Yevgeniy Vorobeychik. 2015. Security Games with Protection Externalities. In *AAAI-15*. Austin, 914–920.
- [14] Chen Greif, Erin Moulding, and Dominique Orban. 2014. Bounds on eigenvalues of matrices arising from interior-point methods. *SIAM Journal on Optimization* 24, 1 (2014), 49–83.
- [15] Alexander Gutfraind, Eric Hagberg, and Feng Pan. 2009. Optimal interdiction of unreactive Markovian evaders. In *CPAIOR-09*. Pittsburgh, 102–116.
- [16] Alexander Gutfraind, Eric A Hagberg, David Izsraelevitz, and Feng Pan. 2011. Interdiction of a Markovian evader. In *Proc. of INFORMS Computing Society*. Monterey, CA.
- [17] Nika Haghtalab, Fei Fang, Thanh Hong Nguyen, Arunesh Sinha, Ariel D Procaccia, and Milind Tambe. 2016. Three Strategies to Success: Learning Adversary Models in Security Games. In *IJCAI-16*. New York, 308–314.
- [18] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NIPS-17*. Long Beach, 1024–1034.
- [19] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. 1983. Stochastic blockmodels: First steps. *Social Networks* 5, 2 (1983), 109–137.
- [20] Manish Jain, Dmytro Korzhuk, Ondřej Vaněk, Vincent Conitzer, Michal Pěchouček, and Milind Tambe. 2011. A double oracle algorithm for zero-sum security games on graphs. In *AAMAS-11*. Taipei, 327–334.
- [21] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR-17*. Toulon.
- [22] Dieter Kraft. 1985. On converting optimal control problems into nonlinear programming problems. In *Computational mathematical programming*. Springer, 261–280.
- [23] Joshua Letchford, Vincent Conitzer, and Kamesh Munagala. 2009. Learning and Approximating the Optimal Strategy to Commit To. In *Algorithmic Game Theory*, Marios Mavronicolas and Vicky G. Papadopoulos (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 250–262.
- [24] Sara Marie Mc Carthy, Milind Tambe, Christopher Kiekintveld, Meredith L Gore, and Alex Killion. 2016. Preventing illegal logging: Simultaneous optimization of resource teams and tactics for security. In *AAAI-16*. New York.
- [25] Richard D McKelvey and Thomas R Palfrey. 1995. Quantal response equilibria for normal form games. *Games and Economic Behavior* 10, 1 (1995), 6–38.
- [26] John Morgan and Felix Vardy. 2004. An experimental study of commitment and observability in Stackelberg games. *Games and Economic Behavior* 49, 2 (2004), 401–423.
- [27] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. 2019. Weisfeiler and Leman go neural: Higher-order graph neural networks. In *AAAI-19*. Honolulu, 4602–4609.
- [28] David P Morton, Feng Pan, and Kevin J Saeger. 2007. Models for nuclear smuggling interdiction. *IEE Transactions* 39, 1 (2007), 3–14.
- [29] Padmanabhan Murugan and Biniam Abebaw. 2014. Factors Contributing to Human Trafficking, Contexts of Vulnerability and Patterns of Victimization: the case of stranded victims in Metema, Ethiopia. *Ethiopian Journal of the Social Sciences and Humanities* 10, 2 (2014), 75–105.
- [30] Michael Victor Nehme. 2009. Two-person games for stochastic network interdiction: models, methods, and complexities. (2009).
- [31] Thanh Hong Nguyen, Rong Yang, Amos Azaria, Sarit Kraus, and Milind Tambe. 2013. Analyzing the Effectiveness of Adversary Modeling in Security Games. In *AAAI-13*. Bellevue, Washington.
- [32] Steven Okamoto, Noam Hazon, and Katia Sycara. 2012. Solving non-zero sum multiagent network flow security games with attack costs. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*. International Foundation for Autonomous Agents and Multiagent Systems, 879–888.
- [33] Feng Pan, William S Charlton, and David P Morton. 2003. A stochastic program for interdicting smuggled nuclear material. In *Network interdiction and stochastic integer programming*. Springer, 1–19.
- [34] Andrei Patrascu and Ion Necoara. 2015. Efficient random coordinate descent algorithms for large-scale structured nonconvex optimization. *Journal of Global Optimization* 61, 1 (2015), 19–46.
- [35] Andrew Perrault, Bryan Wilder, Eric Ewing, Aditya Mate, Bistra Dilkina, and Milind Tambe. 2020. End-to-end Game-focused Learning of Adversary Behavior in Security Games. In *AAAI-2020*. New York.
- [36] Gail Emilia Rosen and Katherine F Smith. 2010. Summarizing the evidence on the international trade in illegal wildlife. *EcoHealth* 7, 1 (2010), 24–32.
- [37] Sankardas Roy, Charles Ellis, Sajjan Shiva, Dipankar Dasgupta, Vivek Shandilya, and Qishi Wu. 2010. A survey of game theory as applied to network security. In *2010 43rd Hawaii International Conference on System Sciences*. IEEE, 1–10.
- [38] Arunesh Sinha, Debarun Kar, and Milind Tambe. 2016. Learning adversary behavior in security games: A PAC model perspective. In *AAMAS-16*. Singapore, 214–222.
- [39] Milind Tambe. 2011. *Security and game theory: algorithms, deployed systems, lessons learned*. Cambridge University Press.
- [40] Iuliana Teodorescu. 2009. Maximum likelihood estimation for Markov Chains. *arXiv preprint arXiv:0905.4131* (2009).
- [41] Paul Tseng. 2001. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications* 109, 3 (2001), 475–494.
- [42] Satya Gautam Vadlamudi, Sailik Sengupta, Marthony Taguinod, Ziming Zhao, Adam Doupe, Gail-Joon Ahn, and Subbarao Kambhampati. 2016. Moving target defense for web applications using bayesian stackelberg games. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 1377–1378.
- [43] Alan Washburn and Kevin Wood. 1995. Two-person zero-sum games for network interdiction. *Operations Research* 43, 2 (1995), 243–251.
- [44] Bernard M Waxman. 1988. Routing of multipoint connections. *IEEE Journal on Selected Areas in Communications* 6, 9 (1988), 1617–1622.
- [45] Rong Yang, Fei Fang, Albert Xin Jiang, Karthik Rajagopal, Milind Tambe, and Rajiv Maheswaran. 2012. Designing better strategies against human adversaries in network security games. In *AAMAS-12*. Valencia.
- [46] Zhengyu Yin, Dmytro Korzhuk, Christopher Kiekintveld, Vincent Conitzer, and Milind Tambe. 2010. Stackelberg vs. Nash in security games: Interchangeability, equivalence, and uniqueness. In *AAMAS-10*. Toronto, 1139–1146.
- [47] Chao Zhang, Arunesh Sinha, and Milind Tambe. 2015. Keeping pace with criminals: Designing patrol allocation against adaptive opportunistic criminals. In *AAMAS-15*. Istanbul, 1351–1359.
- [48] Mara E Zimmerman. 2003. The black market for wildlife: Combating transnational organized crime in the illegal wildlife trade. *Vand. J. Transnat'l L.* 36 (2003), 1657.