

Refinement of Strong Stackelberg Equilibria in Security Games

Bo An, Milind Tambe, Fernando Ordonez, Eric Shieh

University of Southern California
Los Angeles, CA 90089
{boa,tambe,fordon,eshieh}@usc.edu

Christopher Kiekintveld

University of Texas, El Paso
El Paso, TX 79968
cdkiekintveld@utep.edu

Abstract

Given the real-world deployments of attacker-defender Stackelberg security games, robustness to deviations from expected attacker behaviors has now emerged as a critically important issue. This paper provides four key contributions in this context. First, it identifies a fundamentally problematic aspect of current algorithms for security games. It shows that there are many situations where these algorithms face multiple equilibria, and they arbitrarily select one that may hand the defender a significant disadvantage, particularly if the attacker deviates from its equilibrium strategies due to unknown constraints. Second, for important subclasses of security games, it identifies situations where we will face such multiple equilibria. Third, to address these problematic situations, it presents two equilibrium refinement algorithms that can optimize the defender's utility if the attacker deviates from equilibrium strategies. Finally, it experimentally illustrates that the refinement approach achieved significant robustness in consideration of attackers' deviation due to unknown constraints.

Introduction

Game theory is an increasingly important paradigm for reasoning about complex security resource allocation and patrolling problems (Basilico, Gatti, and Amigoni 2009; Korzhyk, Conitzer, and Parr 2010; Dickerson et al. 2010). Much of this work uses Stackelberg game models to represent the commitment that the security forces make to a security policy, and the capability of attackers to use surveillance to learn about the policy during the planning phase of an attack. Two important examples of this are ARMOR used at the LAX Airport to randomize checkpoint placement and canine patrols, and IRIS used by the Federal Air Marshals Service (FAMS) to schedule air marshals (Jain et al. 2010b). Additionally, the United States Transportation Security Administration (TSA) is currently evaluating the GUARDS system for use in scheduling airport security operations (Pita et al. 2011).

To date, the analysis of Stackelberg security games has focused primarily on computing Strong Stackelberg Equilibrium (SSE), and many algorithms have been developed for quickly computing these equilibria in various cases (e.g., DOBSS (Jain et al. 2010b), ERASER (Jain et al.

2010b), ASPEN (Jain et al. 2010a)). However, the assumptions of perfect knowledge and optimal attacker behavior that SSE is based on are very strong, and an important line of recent work focuses on developing more robust solutions for security games. This includes work on both payoff uncertainty in security games (Jain et al. 2010b; Kiekintveld, Tambe, and Marecki 2010) and uncertainty about the behavior of human adversaries (Pita et al. 2010).

In this paper we focus on addressing a kind of uncertainty that has not been studied in the literature, but which may arise in many real-world applications: the possibility that attackers have unknown capability constraints that restrict the set of targets they can feasibly attack. For example, in the airline domain an attacker may not be able to attack certain flights due to an inability to acquire tickets, flight delays/cancellations, visa or passport issues, location constraints, or any number of other issues. The previous work has not developed robust solutions to this form of uncertainty and the arbitrarily selected SSEs may give defenders low utilities if attackers deviate due to unknown constraints.

We develop a solution based on the notion of equilibrium refinement, which is common in the game theory literature (Carlsson and van Damme 1993; Kreps and Wilson 1982; Myerson 1978; Selten 1975). We introduce a new equilibrium refinement for Stackelberg security games based on a dominance criteria motivated by constrained attackers. In cases where there are multiple equilibria, this refinement is able to select more robust equilibria that make the most efficient use of available resources. An important property of the refinement is that it improves robustness to deviations by constrained attackers without any loss in the defender's expected payoff against unconstrained attackers. We also show empirically that the refinement is more robust than standard SSE solutions when there is payoff uncertainty.

We begin by presenting background on Stackelberg security games, and then introduce the key problem of multiple equilibria in security games and define our refinement criteria. We then present a characterization of the cases where multiple equilibria can arise, due to a variety of different restrictions on resources and payoffs. After this theoretical analysis we describe our algorithmic approach for computing the refinement, including specific algorithms for two important cases. We conclude with an experimental evaluation of the refinement algorithms, demonstrating that this tech-

nique significantly improves the robustness of solutions to Stackelberg security games.

Stackelberg Security Games

A generic Stackelberg security game has two players, a defender which first decides how to use m identical resources to protect a set of targets T ($m < |T|$), and an attacker which observes the defender's strategy before choosing a target to attack. A defender's pure strategy is a subset of targets from T such that at most m targets from T are protected. An attacker's pure strategy is a target from T which will be attacked. A mixed strategy allows a player to play a probability distribution over pure strategies. From a mixed strategy of the defender, we can compute the overall coverage of each target. Formally, the defender's mixed strategy can be compactly represented as a coverage vector $\mathbf{c} = \langle c_t \rangle$ where c_t is the probability that target t is covered (Kiekintveld et al. 2009; Yin et al. 2010). The attacker's mixed strategy $\mathbf{a} = \langle a_t \rangle$ is a vector where a_t is the probability of attacking target t .

The payoffs for an agent depend on which target is attacked and how the target is covered. The defender's payoff for an uncovered attack is denoted as $U_d^u(t)$, and for a covered attack $U_d^c(t)$. Similarly, $U_a^u(t)$ and $U_a^c(t)$ are the attacker's payoffs. As a key property of security games, we assume $\Delta_d(t) = U_d^c(t) - U_d^u(t) > 0$ and $\Delta_a(t) = U_a^u(t) - U_a^c(t) > 0$. In other words, adding resources to cover a target hurts the attacker and helps the defender. For a strategy profile $\langle \mathbf{c}, \mathbf{a} \rangle$, the expected utilities for both agents are given by:

$$U_d(\mathbf{c}, \mathbf{a}) = \sum_{t \in T} a_t U_d(\mathbf{c}, t), \text{ where } U_d(\mathbf{c}, t) = c_t U_d^c(t) + (1 - c_t) U_d^u(t)$$

$$U_a(\mathbf{c}, \mathbf{a}) = \sum_{t \in T} a_t U_a(\mathbf{c}, t), \text{ where } U_a(\mathbf{c}, t) = c_t U_a^c(t) + (1 - c_t) U_a^u(t)$$

It thus follows that $U_d^u(t) \leq U_d(\mathbf{c}, t) \leq U_d^c(t)$ and $U_a^c(t) \leq U_a(\mathbf{c}, t) \leq U_a^u(t)$ for any target t .

In a Stackelberg model, the defender chooses its strategy first, and the attacker chooses a strategy after observing the defender's choice. The attacker's response function is $g(\mathbf{c}) : \mathbf{c} \rightarrow \mathbf{a}$. We assume that the $g(\mathbf{c})$ is unique to every \mathbf{c} . The standard solution concept for Stackelberg games is Strong Stackelberg Equilibrium (SSE) (Breton, Alg, and Haurie 1988; Leitmann 1978; von Stengel and Zamir 2004). A pair of strategies $\langle \mathbf{c}, g(\mathbf{c}) \rangle$ form an SSE if they satisfy the following:

1. The defender plays a best-response: $U_d(\mathbf{c}, g(\mathbf{c})) \geq U_d(\mathbf{c}', g(\mathbf{c}'))$ for any \mathbf{c}' .
2. The attacker plays a best-response: $g(\mathbf{c}) \in F_a(\mathbf{c})$ where $F_a(\mathbf{c}) = \arg \max_{\mathbf{a}} U_a(\mathbf{c}, \mathbf{a})$ is the set of follower best-responses.
3. The attacker breaks ties optimally for the defender: $U_d(\mathbf{c}, g(\mathbf{c})) \geq U_d(\mathbf{c}, \mathbf{a}')$ for any $\mathbf{a}' \in F_a(\mathbf{c})$.

There always exists an optimal pure-strategy response for the attacker. Without loss of generality, we restrict to attackers' pure strategies. Given the defender's strategy \mathbf{c} , the *attack set* $\Gamma(\mathbf{c}) = \arg \max_{t \in T} U_a(\mathbf{c}, t)$ contains all targets that

yield the maximum expected payoff for the attacker. Obviously, it follows that $U_a(\mathbf{c}, \mathbf{a}) = U_a(\mathbf{c}, t)$ for any $t \in \Gamma(\mathbf{c})$.

Refinement for SSE in Security Games

A well-known property of SSE is that all SSE give the same expected payoff for the leader (defender) (Breton, Alg, and Haurie 1988; Leitmann 1978), so it may seem strange to be concerned with selecting among the possible SSE solutions. Our interest in a refinement is driven by two observations. First, we find that in many security games based on real-world problems (particularly when there are restrictions on how security resource can be allocated), there are an infinite number of SSE solutions. Worse, in many of these solutions a portion of the available resource are not used productively, since they can be assigned arbitrarily without affecting the solution quality. The second observation is that SSE solutions are not robust to deviations in the strategy of the follower (attacker). This leads to the core idea of our work: using refinement to increase the robustness of SSE solutions, *without* sacrificing solution quality.

To see how multiple equilibria can arise in security games, consider the following example based on the FAMS domain (Jain et al. 2010b). There are 4 targets, representing flights, which are divided into two partitions $\{t_1, t_2\}$ and $\{t_3, t_4\}$, which represent flights leaving from two different airports. There is one air marshal at each airport, so the first marshal can only take flight t_1 or t_2 , and the second can only take t_3 or t_4 . For the payoffs shown below, the game has an infinite number of SSE solutions: $\langle \mathbf{c} = (x, y, 0.5, 0.5), \mathbf{a} = (0, 0, 1, 0) \rangle$ such that $0 \leq x \leq 1$ and $0 \leq y \leq 1 - x$. In effect, it does not matter exactly what the first air marshals does, because the attacker always prefers to attack flight 3 regardless. In contrast, if there was no restriction on the resources, the unique SSE solution is $\langle \mathbf{c} = (1/3, 0, 1, 2/3), \mathbf{a} = (0, 0, 1, 0) \rangle$, which places additional resources on flights 3 and 4.

| | $U_d^u(t)$ | $U_d^c(t)$ | $U_a^u(t)$ | $U_a^c(t)$ |
|-------|------------|------------|------------|------------|
| t_1 | 3 | 4 | 9 | 6 |
| t_2 | 2 | 3 | 7 | 6 |
| t_3 | 4 | 6 | 10 | 8 |
| t_4 | 2 | 3 | 12 | 6 |

It is also possible to have multiple equilibria with no resource restrictions. Consider the game with only targets $\{t_1, t_2, t_3\}$ and 2 resources that can cover any of the three targets. This game also has an infinite number of SSEs $\langle \mathbf{c} = (x, y, 1), \mathbf{a} = (0, 0, 1) \rangle$ such that $x \geq 1/3$ and $y \leq 1 - x$. For instance, $S_1 = \langle \mathbf{c} = (1, 0, 1), \mathbf{a} = (0, 0, 1) \rangle$ and $S_2 = \langle \mathbf{c} = (0.75, 0.25, 1), \mathbf{a} = (0, 0, 1) \rangle$. In every case, the attacker chooses to attack the fully covered target t_3 with a utility of 8. The intuition for this case is that the attacker prefers to attack a highly valuable target, even though it is heavily defended.¹ This is an important real-world phenomenon. Even though airports, government buildings, and

¹We note that the security game model as presented focuses on allocating a single type of resource (e.g., air marshals or checkpoints). The addition of altogether different types of protection could further enhance security; our model does not address this direction, but does provide useful diagnostic information to identify cases where this may be useful.

other critical targets are very heavily secured, they may still be more desirable targets than millions of unsecured “soft targets”. In practice, it is never possible to achieve perfect security, so there is always some chance of success—and even a failed attack against an important target may result in a costly response and a large amount of publicity.

We will further characterize the situations where multiple SSE are possible in the next section. Now, we turn to the question of how to select among the multiple equilibria. Our primary motivation is to increase the robustness of the solution in case the attacker deviates from the SSE strategy. This could occur for many different reasons, including bounded rationality, payoff noise, and unknown capability constraints. We formulate our refinement based on the assumption that the attacker deviates due to unknown constraints for three reasons: (1) it has not been previously studied in the literature so there are no known methods for dealing with this form of uncertainty, (2) it leads to a natural refinement criteria, and (3) it is real-world as discussed at the beginning. In our experimental results, we show that the refinement is also effective for deviations as a result of payoff uncertainty.

Consider again the example above with no resource restrictions. All SSEs give the defender the same optimal utility 6, but they give dramatically different results if the attacker deviates from the equilibrium strategy. Suppose that the attacker is not able to attack t_3 for some unknown reason. In S_1 , the attacker’s second-best response is t_2 and the defender’s utility is 2. For S_2 , the attacker is indifferent between t_1 and t_2 . Since it breaks ties optimally for the defender, it attacks t_1 and the defender’s utility is 3.75. In this example S_2 is more robust than S_1 , and this robustness is “free” in the sense that the defender still receives the optimal utility if the attacker attacks its optimal target t_3 .

More formally, we consider a model of a *constrained attacker*. The constrained attacker has the same payoff function as defined in the base security game, but cannot attack a subset of the possible targets. Inspired by refinement solution concepts that depend on small probability of deviations (e.g., (Selten 1975)), we assume that each target appearing in the set of targets cannot be attacked with small probability, which could be different for different targets. This simple model has several desirable properties. It is unlikely that the attacker will be unable to attack any particular target, and increasingly unlikely that the attacker will be unable to attack larger combinations of specific targets. Moreover, even when the attacker is forced to deviate due to a constraint, he still behaves intelligently and chooses the next-best alternative rather than acting randomly.

Based on this model, we now define our equilibrium refinement concept. Given an SSE $\langle \mathbf{c}, \mathbf{a} \rangle$, we define an ordering over the targets as follows. Let target $t(1)$ be the target that will be attacked if the attacker is unconstrained. Let target $t(i)$ be the target that will be attacked if the attacker is not able to attack targets $t(1), \dots, t(i-1)$. Utility vector $\mathbf{v} = \langle v_i \rangle$ represents the defender’s utilities where v_i is the defender’s utility if target $t(i)$ is attacked, i.e., $v_i = c_{t(i)} U_d^c(t(i)) + (1 - c_{t(i)}) U_d^u(t(i))$. We define a dominance relation between SSEs based on their utility vectors.

Definition 1. Given two SSEs $\langle \mathbf{c}, \mathbf{a} \rangle$ and $\langle \mathbf{c}', \mathbf{a}' \rangle$, we have two utility vectors \mathbf{v} and \mathbf{v}' . SSE $\langle \mathbf{c}, \mathbf{a} \rangle$ **dominates** SSE $\langle \mathbf{c}', \mathbf{a}' \rangle$ if there exists i such that 1) $v_i > v'_i$ and 2) for all $1 \leq j < i$, $v_j = v'_j$.

The refinement criterion is to find an SSE that is not dominated by any other SSE. The criteria demands that the solution first optimizes for the most likely scenario, in which the attacker can attack the optimal SSE target. The secondary criteria is to optimize against an attacker that chooses the second-best target, which is the next most-likely scenario. Each successive scenario is less likely, but the refinement continues to select an optimal response as long as resources are available.

Characterization of Multiple SSEs

We wish to characterize situations where multiple SSEs arise, and to that end begin with analyzing the relationship between the coverage vector and the attack set in an SSE. Our analysis assumes no maximum resource restrictions and we show how to generalize the analysis to games with maximum resource restrictions at the end of this section.

Lemma 2. Assume that there is an SSE $\langle \mathbf{c}, \mathbf{a} \rangle$. If the coverage of each target $t \in \Gamma(\mathbf{c})$ is less than 1, i.e., $c_t < 1$, then $c_{t'} = 0$ for any target $t' \notin \Gamma(\mathbf{c})$. In addition, $\sum_{t \in \Gamma(\mathbf{c})} c_t = m$.

Proof. Assume that $c_{t'} > 0$ for a $t' \notin \Gamma(\mathbf{c})$. Since $U_a(\mathbf{c}, \mathbf{a}) > U_a(\mathbf{c}, t')$, we can slightly reduce the coverage of target t' and increase the coverage of targets $\Gamma(\mathbf{c})$ such that the attacker is still indifferent among the attack set $\Gamma(\mathbf{c})$ and no new target is added to $\Gamma(\mathbf{c})$. Then the defender’s utility increases, which contradicts to the fact that $\langle \mathbf{c}, \mathbf{a} \rangle$ is an SSE. Therefore, $c_{t'} = 0$ for any target $t' \notin \Gamma(\mathbf{c})$. The same argument shows that $\sum_{t \in \Gamma(\mathbf{c})} c_t = m$ if $c_t < 1$ for every $t \in \Gamma(\mathbf{c})$. \square

Lemma 2 implies that if a target $t' \notin \Gamma(\mathbf{c})$ is covered, i.e., $c_{t'} > 0$, at least one target $t \in \Gamma(\mathbf{c})$ is fully covered, i.e., $c_t = 1$.

Lemma 3. Assume that there are two targets t and t' such that $U_a^u(t) \geq U_a^u(t')$ and there is an SSE $\langle \mathbf{c}, \mathbf{a} \rangle$. If $t' \in \Gamma(\mathbf{c})$ and $t \notin \Gamma(\mathbf{c})$, there is at least one target $t'' \in \Gamma(\mathbf{c})$ such that $c_{t''} = 1$ and there is another SSE $\langle \mathbf{c}', \mathbf{a} \rangle$ in which $t \in \Gamma(\mathbf{c}')$.

Proof. Since target t' is in the attack set $\Gamma(\mathbf{c})$, we have 1) $U_a(\mathbf{c}, \mathbf{a}) = U_a(\mathbf{c}, t')$ according to the definition of attack sets and 2) $U_a(\mathbf{c}, t') \leq U_a^u(t')$ since $c_{t'} \geq 0$. Since $U_a^u(t) \geq U_a^u(t')$, it follows that $U_a(\mathbf{c}, \mathbf{a}) \leq U_a^u(t)$. If $c_t = 0$, the target t should be in the attack set since $U_a(\mathbf{c}, \mathbf{a}) \leq U_a^u(t) = U_a(\mathbf{c}, t)$. Therefore, $c_t > 0$. It then follows that at least one target $t'' \in \Gamma(\mathbf{c})$ is fully covered, i.e., $c_{t''} = 1$ (Lemma 2).

Since $t \notin \Gamma(\mathbf{c})$, it follows that $U_a(\mathbf{c}, \mathbf{a}) > U_a(\mathbf{c}, t)$, i.e., $c_t > (U_a^u(t) - U_a(\mathbf{c}, \mathbf{a})) / \Delta_a(t)$. If we reduce the coverage c_t to $c'_t = (U_a^u(t) - U_a(\mathbf{c}, \mathbf{a})) / \Delta_a(t)$, we construct another SSE $\langle \mathbf{c}', \mathbf{a} \rangle$ in which $\mathbf{c}' = (\mathbf{c}_{-t}, c'_t)$ where \mathbf{c}_{-t} is the coverage of all targets except for target t in the SSE $\langle \mathbf{c}, \mathbf{a} \rangle$. \square

Lemma 3 implies that we may create another SSE with a larger attack set by reducing coverage of targets outside the attack set. The attack set of SSE S_1 in our example is $\{t_3\}$. If we reduce the coverage of t_1 to $1/3$, we have another SSE $\langle (1/3, 0, 1), (0, 0, 1) \rangle$ with a larger attack set $\{t_1, t_3\}$. We

Algorithm 1: Compute the maximum attack set for a security game with m resources and targets T

```

1  $i \leftarrow 0, \mathcal{M} \leftarrow \emptyset, T$  are sorted by  $U_a^u(t)$ ;
2 while  $i \leq |T|$  do
3   if  $\mathcal{M} = T$  then return  $\mathcal{M}$ ;
4    $j \leftarrow i + 1, \mathcal{M}' \leftarrow \mathcal{M} \cup \{t_j\}$ ;
5   while  $j < |T|$  and  $U_a^u(t_{j+1}) = U_a^u(t_j)$  do
6      $\mathcal{M}' \leftarrow \mathcal{M}' \cup \{t_{j+1}\}; j \leftarrow j + 1$ ;
7   end
8   if Condition C1 or C2 is violated for attack set  $\mathcal{M}'$  then
9     return  $\mathcal{M}$ ;
9    $\mathcal{M} \leftarrow \mathcal{M}'; i \leftarrow j$ ;
10 end

```

define *maximum attack set SSE* (MSSE). An SSE $\langle \mathbf{c}, \mathbf{a} \rangle$ is an MSSE if $\Gamma(\mathbf{c}) = \{t : U_a^u(t) \geq U_a(\mathbf{c}, \mathbf{a})\}$. Without loss of generality, we assume that targets $T = \{t_1, \dots, t_{|T|}\}$ are sorted by $U_a^u(t)$, i.e., $U_a^u(t_i) \geq U_a^u(t_j)$ if $1 \leq i \leq j \leq |T|$.

Theorem 4. Any security game could not have two MSSEs with different attack sets.

Proof. Assume that there are two MSSEs $\langle \mathbf{c}, \mathbf{a} \rangle$ and $\langle \mathbf{c}', \mathbf{a}' \rangle$ such that $\Gamma(\mathbf{c}) \neq \Gamma(\mathbf{c}')$. Given the definition of MSSE, either $\Gamma(\mathbf{c}) \subset \Gamma(\mathbf{c}')$ or $\Gamma(\mathbf{c}') \subset \Gamma(\mathbf{c})$ should be satisfied. We assume that $\Gamma(\mathbf{c}) \subset \Gamma(\mathbf{c}')$ and the other situation can be proved in the same way. Let $x = U_a(\mathbf{c}, \mathbf{a})$ and $x' = U_a(\mathbf{c}', \mathbf{a}')$.

Assume that $x > x'$. Given that $x = U_a^u(t) - \Delta_a(t)c_t$ and $x' = U_a^u(t) - \Delta_a(t)c'_t$, we have $c_t < c'_t$ for all $t \in \Gamma(\mathbf{c})$. Therefore, $c_t < 1$ for all $t \in \Gamma(\mathbf{c})$. By Lemma 2, this means $m = \sum_{t \in \Gamma(\mathbf{c})} c_t < \sum_{t \in \Gamma(\mathbf{c}')} c'_t \leq m$, which is a contradiction.

Likewise, we can show a contradiction if $x < x'$. Thus $x' = x$ and MSSEs have the same maximum attack set. \square

Algorithm 1 computes the unique maximum attack set for all MSSEs. The algorithm starts with the attack set $\mathcal{M} = \{t_1\}$ and gradually adds targets into \mathcal{M} . Given a maximum attack set \mathcal{M} for an MSSE $\langle \mathbf{c}, \mathbf{a} \rangle$, it follows that $U_a(\mathbf{c}, \mathbf{a}) \leq \min_{t \in \mathcal{M}} U_a^u(t)$. Accordingly, the minimum coverage of each target $t \in \mathcal{M}$ is $(U_a^u(t) - x)/\Delta_a(t)$ where $x = \min_{t \in \mathcal{M}} U_a^u(t)$. Therefore, the minimum number of resources needed is $\sum_{t \in \mathcal{M}} (U_a^u(t) - x)/\Delta_a(t)$. By adding a target t_i to the attack set \mathcal{M} , the coverage of each target $t \in \mathcal{M}$ and the number of resources needed will not decrease. In particular, the defender's utility will not decrease since the attacker will choose the target best for the defender from the attack set. Clearly, the attack set of each MSSE is the largest attack set \mathcal{M} satisfying conditions **C1** and **C2** ($c_t = (U_a^u(t) - \min_{t \in \mathcal{M}} U_a^u(t))/\Delta_a(t)$ for $t \in \mathcal{M}$):

- **C1:** $\sum_{t \in \mathcal{M}} c_t \leq m$
- **C2:** $c_t \leq 1$ for each $t \in \mathcal{M}$

The next theorem gives the necessary and sufficient conditions for the uniqueness of SSEs.

Theorem 5. A security game has a unique SSE if and only if (\mathcal{M} is the maximum attack set computed using Algorithm 1)

$$\max_{t \in \mathcal{M}} U_a^c(t) \leq \frac{\sum_{t \in \mathcal{M}} \frac{U_a^u(t)}{\Delta_a(t)} - m}{\sum_{t \in \mathcal{M}} \frac{1}{\Delta_a(t)}}$$

Proof. Given a maximum attack set \mathcal{M} , we can compute all SSEs as follows. Since decreasing the attacker's utility in attacking \mathcal{M} implies increasing the defender's utility, the defender will use as many resources as possible to cover \mathcal{M} while satisfying conditions **C1** and **C2**. Assume that the attacker's lowest utility for attacking any target in \mathcal{M} is x . Then the coverage of each target $t \in \mathcal{M}$ is $c_t = (U_a^u(t) - x)/\Delta_a(t)$. If we only consider **C1** ($\sum_{t \in \mathcal{M}} c_t \leq m$), we have

$$x \geq \underline{x} = \frac{\sum_{t \in \mathcal{M}} \frac{U_a^u(t)}{\Delta_a(t)} - m}{\sum_{t \in \mathcal{M}} \frac{1}{\Delta_a(t)}}$$

Condition **C2** ($c_t \leq 1$ for each $t \in \mathcal{M}$) implies that $x \geq U_a^c(t)$ for each $t \in \mathcal{M}$, i.e., $x \geq \bar{x} = \max_{t \in \mathcal{M}} U_a^c(t)$. Therefore, the attacker's utility in attacking \mathcal{M} is $x = \max\{\bar{x}, \underline{x}\}$. The optimal coverage of each target $t \in \mathcal{M}$ is $c_t = (U_a^u(t) - \max\{\bar{x}, \underline{x}\})/\Delta_a(t)$. If $\bar{x} \leq \underline{x}$, all resources will be used and there is a unique SSE, which is also an MSSE. \square

If $\bar{x} \leq \underline{x}$, the coverage vector \mathbf{c} for the unique SSE is: $c_t = \frac{U_a^u(t) - \bar{x}}{\Delta_a(t)}$ if $t \in \mathcal{M}$, and $c_t = 0$ otherwise. In the unique (M)SSE, the attacker will attack the target $t^* \in \mathcal{M}$ which gives the defender the highest utility.

If $\bar{x} > \underline{x}$, the attacker's expected utility of attacking any target in an SSE is \bar{x} . In an MSSE, the coverage of each target $t \in \mathcal{M}$ is $c_t = (U_a^u(t) - \bar{x})/\Delta_a(t)$. Given the Algorithm 1, it follows that $m - \sum_{t \in \mathcal{M}} c_t > 0$. The remaining resources $m - \sum_{t \in \mathcal{M}} c_t$ could be arbitrarily allocated to targets $T \setminus \mathcal{M}$.

If $\bar{x} > \underline{x}$, there are also an infinite number of SSEs with attack sets which are subsets of \mathcal{M} . Since the attacker will attack the target t^* which can give the defender the highest utility, t^* should be in the attack set of for any SSE, i.e., the coverage of t^* is $(U_a^u(t^*) - \bar{x})/\Delta_a(t^*)$. For $t \in \mathcal{M} \setminus \{t^*\}$, the coverage c_t should be no less than $(U_a^u(t) - \bar{x})/\Delta_a(t)$. Formally, the coverage vectors for all SSEs are given as follows.

$$c_t \begin{cases} = \frac{U_a^u(t) - \bar{x}}{\Delta_a(t)} & \text{if } t = t^* \\ \in [\frac{U_a^u(t) - \bar{x}}{\Delta_a(t)}, 1] \text{ s.t. } \sum_{t' \in \mathcal{M}} c_{t'} \leq m & \text{if } t \in \mathcal{M} \setminus \{t^*\} \\ \in [0, 1] \text{ s.t. } \sum_{t' \notin \mathcal{M}} c_{t'} \leq m - \sum_{t' \in \mathcal{M}} c_{t'} & \text{if } t \notin \mathcal{M} \end{cases}$$

The above analysis assumes no resource restrictions and cannot be generalized to games with *arbitrary* scheduling restrictions Θ , e.g., minimum/maximum coverage requirement for some targets. With arbitrary restrictions, it's impossible to characterize SSE uniqueness since restrictions themselves (e.g., $c_t = 0.5$ for each target t) could uniquely define the defender's strategy. However, the above analysis can be applied to security games with maximum resource restrictions as in the FAMS domain. Each restriction θ provides an upper bound w_θ on the total coverage of a set of targets T_θ , i.e., $\sum_{t \in T_\theta} c_t \leq w_\theta$. When we use Algorithm 1 to compute the maximum attack set \mathcal{M} , we need to also guarantee that condition **C3** is not violated for \mathcal{M} :

- **C3:** $\sum_{t \in T_\theta} c_t \leq w_\theta$ for each θ

where $c_t = (U_a^u(t) - \min_{t \in \mathcal{M}} U_a^u(t))/\Delta_a(t)$ for $t \in \mathcal{M}$ and $c_t = 0$ for $t \notin \mathcal{M}$. The attacker's utility while attacking

Algorithm 2: Compute a refined SSE for a game with m resources and targets T

```

1  $i \leftarrow 0, m' \leftarrow 0;$ 
2 while true do
3   Compute an SSE  $\langle \mathbf{c}, \mathbf{a} \rangle$  using ERASER;
4   Let  $x \leftarrow c_{t(i)}$  if  $a_{t(i)} = 1$ ;
5   Add restrictions  $a_{t(i)} = 0$  and  $c_{t(i)} = x$  to  $\Theta$ ;
6   Add restrictions that the attacker's expected utility for any
   target must be no higher than the expected payoff for  $t(i)$ ;
7    $i ++, m' \leftarrow m' + x$ ;
8   if  $i = |T|$  or  $m' = m$  then return  $\langle \mathbf{c}, g(\mathbf{c}) \rangle$ ;
9 end

```

the maximum attack set \mathcal{M} is the lowest utility such that conditions C1 – C3 are satisfied. Then we can compute the coverage vector for targets \mathcal{M} . There are multiple SSEs if there are still resources left after covering \mathcal{M} .

Algorithms for Refinement

We now describe a general framework for computing the equilibrium refinement and two instantiations of this framework for specific classes of security games. As noted previously there are many existing algorithms for finding SSE in various types of security games. We can compute the refined solution by iteratively applying the SSE solver to restricted forms of the security game. Algorithm 2 gives the pseudocode for this approach using the ERASER MILP (Jain et al. 2010b) as the base SSE solver. The algorithm first computes an SSE $\langle \mathbf{c}, \mathbf{a} \rangle$ for the unmodified game with m resources and targets T . Let $t(1)$ be the target attacked in this SSE, with a coverage of $c_{t(1)} = x$. For the second iteration, we assume that the attacker will not attack $t(1)$, but still maintain the equilibrium conditions from the original solutions. Specifically, we add to the original optimization problem constraints that the coverage $c_{t(1)} = x$, that the attacker does not attack $t(1)$, and that for all targets other than $t(1)$, the attacker's expected utility is less than or equal to the expected utility for attacking $t(1)$. We solve the modified MILP to compute a new SSE $\langle \mathbf{c}', \mathbf{a}' \rangle$ where the attacker chooses a new target $t(2)$, and a new set of constraints is added for the third iteration. This process continues until there are no targets left or there are no unconstrained resources in the restricted problem on some iteration; the maximum number of iterations is $|T|$.

Theorem 6. *The strategy profile computed by Algorithm 2 is an SSE and is not dominated by another SSE.*

Proof. We first show that the strategy profile $\langle \mathbf{c}, \mathbf{a} \rangle$ returned by Algorithm 2 is an SSE. Assume that the algorithm stops after n iterations and at iteration i , the target will be attacked is k_i . Assume that the SSE computed in round 1 is $\langle \mathbf{c}', \mathbf{a}' \rangle$, i.e., $a'_{k_1} = 1$. By construction, the attacker's utility by attacking target k_1 is no less than that by attacking targets k_2, \dots, k_n , which implies that $k_1 \in \Gamma(\mathbf{c})$. This, combined with the fact that (since $a'_{k_1} = 1$) $U_d(\mathbf{c}, \mathbf{a}) = U_d(\mathbf{c}', k_1) = U_d(\mathbf{c}', \mathbf{a}')$, implies that $\langle \mathbf{c}, \mathbf{a} \rangle$ is an SSE.

Algorithm 3: $\mathbf{RSSE}(\Upsilon(m, T))$

```

1 Compute the maximum attack set  $\mathcal{M}$  using Algorithm 1;
2 Let  $\langle \mathbf{c}, \mathbf{a} \rangle$  be an MSSE (Theorem 5);
3 if There is no other SSE other than  $\langle \mathbf{c}, \mathbf{a} \rangle$  (Theorem 5) then
4   return  $\langle \mathbf{c}, \mathbf{a} \rangle$ ;
5 else
6   Let  $t(i)$  be the target such that  $a_{t(i)} = 1$ ;
7   Set  $T' = T \setminus \{t(i)\}$  and  $m' = m - \sum_{t \in T \setminus T'} c_t$ ;
8   Let  $\langle \mathbf{c}', \mathbf{a}' \rangle = \mathbf{RSSE}(\Upsilon(m', T'))$ ;
9   return  $\langle \mathbf{c}_{T \setminus T'} \cup \mathbf{c}', \mathbf{a} \rangle$ ;
10 end

```

The non-dominance of the SSE $\langle \mathbf{c}, \mathbf{a} \rangle$ is guaranteed by construction. We first compute the defender's optimal coverage. Then we fix the coverage for the target that will be attacked and restrict the attacker's utility in attacking other targets. Thus the defender can gain the highest utility if the attacker is attacking its best target. Then we compute the defender's best strategy in the remaining game, which can give the attacker the highest utility if the attacker attacks the second best target. This process continues until the remaining game has only one SSE. \square

This general iterative algorithm structure can be used with other SSE solvers as well. For example, Algorithm 3 computes a refined SSE for games where the defender resources have no restrictions. This time we use Algorithm 1 as the base SSE solver, which is much faster than ERASER for this restricted class of games. For a game $\Upsilon(m, T)$, we first compute the maximum attack set \mathcal{M} . If there is a unique SSE $\langle \mathbf{c}, \mathbf{a} \rangle$, the algorithm terminates and returns $\langle \mathbf{c}, \mathbf{a} \rangle$. Otherwise, it will compute an MSSE $\langle \mathbf{c}, \mathbf{a} \rangle$ in which the attacker will attack target $t(i)$ in iteration i . We fix the coverage of $t(i)$ and compute the refined SSE $\langle \mathbf{c}', \mathbf{a}' \rangle$ for the remaining game with targets $T \setminus \{t(i)\}$ and $m - c_{t(i)}$ resources. In the final refined SSE the coverage of target t is c_t and the coverage for targets $T \setminus \{t\}$ other than t is \mathbf{c}' .

Theorem 7. *The strategy profile computed by Algorithm 3 is an SSE and is not dominated by any other SSE.*

Proof. We first show that the strategy profile returned by Algorithm 3 is an SSE, which is obviously true if the game $\Upsilon(m, T)$ has a unique SSE. If $\Upsilon(m, T)$ has multiple SSEs and $\langle \mathbf{c}, \mathbf{a} \rangle$ is an MSSE with $a_t = 1$, showing that the combined strategy profile $\langle \mathbf{c}_t \cup \mathbf{c}', \mathbf{a} \rangle$ is an SSE is reduced to showing that $\mathbf{a} \in F_a(\mathbf{c}_t \cup \mathbf{c}')$. This, combined with the fact that (since $a_t = 1$) $U_d(\mathbf{c}_t \cup \mathbf{c}', \mathbf{a}) = U_d(\mathbf{c}, t) = U_d(\mathbf{c}, \mathbf{a})$, implies that $\langle \mathbf{c}_t \cup \mathbf{c}', \mathbf{a} \rangle$ is an SSE.

Note that the set $\mathcal{M} \setminus \{t\}$ is a set of targets with the largest uncovered attacker utilities and they can form an attack set (just set $\mathbf{c}' = \mathbf{c}_{-t}$). This means that the set $\mathcal{M} \setminus \{t\}$ would be constructed during the execution of Algorithm 1 on $\mathbf{RSSE}(\Upsilon(m - c_t, T \setminus \{t\}))$ and therefore $\mathcal{M} \setminus \{t\} \subseteq \mathcal{M}'$ where \mathcal{M}' be the maximum attack set for the remaining game $\Upsilon(m - c_t, T \setminus \{t\})$. In the remainder, Algorithm 1 decreases x which increases the coverage on the targets $\mathcal{M} \setminus \{t\}$. Therefore we have that $c'_{t'} \geq c_{t'}$ for all $t' \in \mathcal{M} \setminus \{t\}$. This implies that $U_a(\mathbf{c}_t \cup \mathbf{c}', t') = U_a^u(t') - c'_{t'} \Delta_a(t') \leq$

$U_a^u(t') - c_{t'} \Delta_a(t') = U_a(\mathbf{c}, t') = U_a(\mathbf{c}, t) = U_a(\mathbf{c}_t \cup \mathbf{c}', t)$. This, combined with the fact that $U_a(\mathbf{c}, \mathbf{a}) > U_a^u(t'')$ for any target $t'' \in T \setminus \mathcal{M}$, shows that \mathbf{a} , with $a_t = 1$, satisfies $\mathbf{a} \in F_a(\mathbf{c}_t \cup \mathbf{c}')$.

The non-dominance of the SSE $\langle \mathbf{c}_t \cup \mathbf{c}', \mathbf{a} \rangle$ is guaranteed by construction. We first compute the defender's optimal strategy. Then we remove the attacker's optimal target and fix the coverage for that target. By doing so, the defender can gain the highest utility if the attacker is attacking its best target. Then we compute the defender's best strategy in the remaining game, which gives the attacker the highest utility if the attacker is attacking the second best target. This process continues until the remaining game has only one SSE. \square

We also note that the refined SSE computed by Algorithm 3 is unique, since there is a unique optimal target (and coverage) at each iteration, leading to a unique restricted game for the next iteration. To find the exact refinement, we need to solve at most $|T|$ games since after one iteration, one target's coverage will be fixed. We can speed up the algorithm further using the following observation. Let $t \in \mathcal{M}$ be the target that will be attacked in all MSSEs. By Theorem 5, there should be a target $t' \in \mathcal{M}$ such that $c_{t'} = 1$. If $t \neq t'$, it is easy to see that in the equilibrium $\langle \mathbf{c}', \mathbf{a}' \rangle$ for $\Upsilon(m - c_t, T \setminus \{t\})$, we have 1) $c'_{t'} = 1$, 2) $c'_{t''} = c_{t''}$ for any target $t'' \in \mathcal{M} \setminus \{t\}$ and 3) $\mathcal{M} \setminus \{t\} = \mathcal{M}'$ where \mathcal{M}' is the maximum attack set of equilibrium $\langle \mathbf{c}', \mathbf{a}' \rangle$. The reason is that the coverage of t' cannot be increased any more. Clearly, the attacker's optimal target t'' in equilibrium $\langle \mathbf{c}', \mathbf{a}' \rangle$ is $t'' = \arg \max_{t'' \in \mathcal{M} \setminus \{t\}} U_d(\mathbf{c}, t'')$. By induction, $c'_{t''} = c_{t''}$ for any target $t'' \in \mathcal{M}$ such that $U_d(\mathbf{c}, t'') \geq \min_{t' \in \mathcal{M}, c_{t'}=1} U_d(\mathbf{c}, t')$. Then in Algorithm 3, we can set $T' = T \setminus \Psi$ where $\Psi = \{t'' : U_d(\mathbf{c}, t'') \geq \min_{t' \in \mathcal{M}, c_{t'}=1} U_d(\mathbf{c}, t'), t'' \in \mathcal{M}\}$. Targets with full coverage are always removed, so we need to solve at most $\min\{m, |T|\}$ games. We used this approach to speed up Algorithm 3 in the experiments.

Experimental Evaluation

We run experiments to test the robustness of refined SSE and "standard" SSE selected arbitrarily by solving the ERASER MILP with default CPLEX settings. Our experiments use 100 sample game instances with 5 defender resources, varying numbers of targets, and randomly-generated payoffs. All payoffs are in $[0, 100]$, and we enforce the constraint that rewards are higher than penalties when the uncovered utility is drawn for each player. For games with restrictions on defender resources, we randomly generate 2 min and 2 max coverage restrictions for sets of 2-4 randomly-chosen targets. The max restriction is set to $1.5n * \frac{m}{|T|}$ where n is the number of targets in the restriction set, and the min restriction is set randomly $U[0.1, 0.2]$. T-test yields p-value < 0.0001 for all comparison of refined SSEs against random SSEs.

Figures 1(a) and 2(a) compare the expected defender utilities for the first 5 targets in the dominance ordering for standard SSE and refined SSE computed by Algorithms 2 and 3, respectively. The results show that (1) refined SSE and standard SSE give the same utility when attackers choose the best target, and (2) that refined SSE gives a much higher defender utility when attackers choose $t(2) - t(5)$.

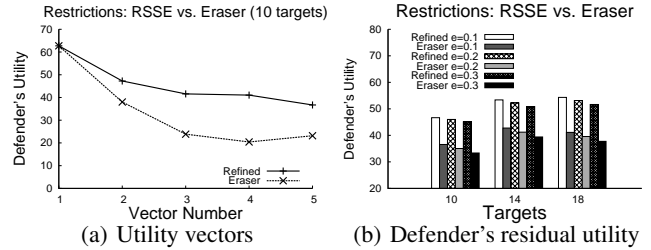


Figure 1: Robustness in games with restrictions.

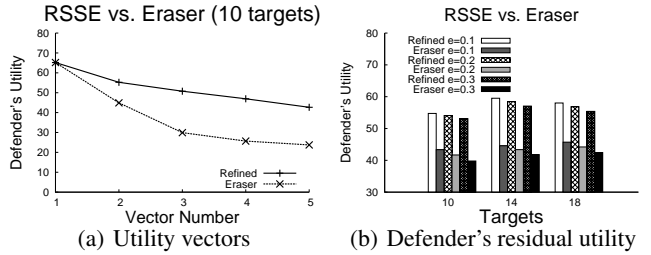


Figure 2: Robustness in games without restrictions.

We can compute the effects of this additional robustness for the constrained attacker model, assuming that the attacker cannot attack any particular target with probability e . Since the values for $t(1)$ are identical we present the *residual* expected utility for the remaining targets. Assuming the attacker does not attack $t(1)$, it will attack $t(2)$ with probability $1 - e$ and attack $t(3)$ with probability $e(1 - e)$. Given the utility vector \mathbf{v} for an SSE, the defender's *residual* utility is given by $\sum_{2 \leq i \leq |T|} (1 - e)e^{i-2} v_i$. Figures. 2(b) and 1(b) compare the residual utilities for refined SSE and standard SSE. The residual utilities for the refined SSE are much higher than standard SSE for any number of targets in our experiments, and the effect increases as the probability e of attacker constraints increases.

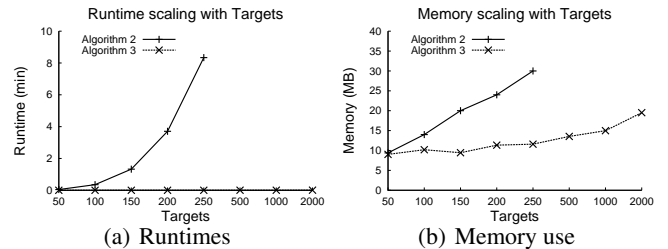


Figure 3: Runtime and memory scaling.

Fig. 3(a) compares runtime performance of Algorithms 2 and 3 for games without defender resource restrictions. The x-axis is the size of the game (in targets), and the y-axis is runtime in minutes. Algorithm 2 has an average runtime of 8m on problems with 250 targets but Algorithm 3 only needs 0.2s. Fig. 3(b) compares the memory performance on the same set of games and shows a similar trend.

We also tested the performance of the refinement to attacker deviations due to payoff uncertainty (rather than attacker constraints). We added mean-0 Gaussian noise to each attacker payoff, with standard deviations chosen randomly from either $U[0, 0.5]$ or $U[0.2, 1.5]$ to generate class-

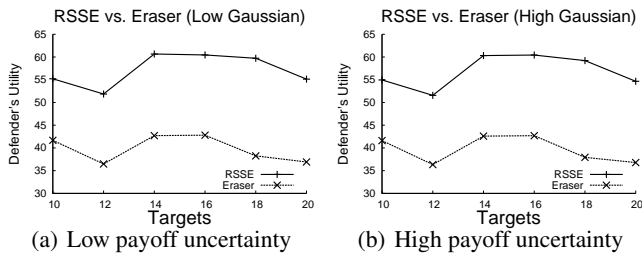


Figure 4: Robustness to payoff uncertainty.

es with “low” or “high” uncertainty. The attacker chooses an optimal response based on the perturbed payoffs. We sample the noisy payoffs and choose a response 1000 times to get a distribution of possible responses for the attacker. Fig. 4 shows the defender’s expected payoffs in this setting for the same class of games described above with no defender resource restrictions. Even though the refinement is designed primarily to optimize against constrained attackers, our results show that it provides additional robustness against deviations due to payoff uncertainty as well, while still providing an optimal payoff if there is no uncertainty.

Conclusions

Robustness is a critically important issue in attacker-defender Stackelberg security games, given that in real-world deployments, attackers may deviate from expected behaviors. This paper provides four key contributions towards such robustness: (1) it identifies a fundamental problem in current security game algorithms, showing that they may select an arbitrary SSE from multiple possible equilibria, handing the defender a significant disadvantage if the attacker deviates from its expected best response due to unknown capability constraints; (2) for important subclasses of security games, it mathematically characterizes situations where the defender will face such problematic multiple equilibria; (3) to address these problematic situations, this paper presents two equilibrium refinement algorithms that optimize defender utility if the attacker deviates from equilibrium strategies; (4) it experimentally illustrates that the refinement approach achieved significant robustness when attackers deviate due to unknown capability constraints, providing a more productive use of available resources.

Acknowledgments

This research was supported under grant 2009-ST-061-CCI002-02 from the Department of Homeland Security to Rutgers University.

References

Basilico, N.; Gatti, N.; and Amigoni, F. 2009. Leader-follower strategies for robotic patrolling in environments with arbitrary topologies. In *Proc. of The 8th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 500–503.

Breton, M.; Alg, A.; and Haurie, A. 1988. Sequential Stackelberg equilibria in two-person games. *Optimization Theory and Applications* 59(1):71–97.

Carlsson, H., and van Damme, E. 1993. Global games and equilibrium selection. *Econometrica* 61(5):989–1018.

Dickerson, J. P.; Simari, G. I.; Subrahmanian, V. S.; and Kraus, S. 2010. A graph-theoretic approach to protect static and moving targets from adversaries. In *Proc. of The 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 299–306.

Jain, M.; Kardes, E.; Kiekintveld, C.; Ordonez, F.; and Tambe, M. 2010a. Security games with arbitrary schedules: A branch and price approach. In *Proc. of The 24th AAAI Conference on Artificial Intelligence*, 792–797.

Jain, M.; Tsai, J.; Pita, J.; Kiekintveld, C.; Rathi, S.; Tambe, M.; and Ordonez, F. 2010b. Software assistants for randomized patrol planning for the LAX airport police and the federal air marshal service. *Interfaces* 40:267–290.

Kiekintveld, C.; Jain, M.; Tsai, J.; Pita, J.; Tambe, M.; and Ordonez, F. 2009. Computing optimal randomized resource allocations for massive security games. In *Proc. of The 8th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 689–696.

Kiekintveld, C.; Tambe, M.; and Marecki, J. 2010. Robust Bayesian methods for Stackelberg security games. In *Proc. of the 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 1467–1468.

Korzhyk, D.; Conitzer, V.; and Parr, R. 2010. Complexity of computing optimal Stackelberg strategies in security resource allocation games. In *Proc. of The 24th AAAI Conference on Artificial Intelligence*, 805–810.

Kreps, D., and Wilson, R. 1982. Sequential equilibria. *Econometrica* 50(4):863–894.

Leitmann, G. 1978. On generalized Stackelberg strategies. *Optimization Theory and Applications* 26(4):637–643.

Myerson, R. 1978. Refinements of the Nash equilibrium concept. *International Journal of Game Theory* 15:133–154.

Pita, J.; Jain, M.; Tambe, M.; Ordóñez, F.; and Kraus, S. 2010. Robust solutions to Stackelberg games: Addressing bounded rationality and limited observations in human cognition. *Artificial Intelligence* 174(15):1142–1171.

Pita, J.; Tambe, M.; Kiekintveld, C.; Cullen, S.; and Steigerwald, E. 2011. GUARDS - game theoretic security allocation on a national scale. In *Proc. of The 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*.

Selten, R. 1975. A reexamination of the perfectness concept for equilibrium points in extensive games. *International Journal of Game Theory* 4:25–55.

von Stengel, B., and Zamir, S. 2004. Leadership with commitment to mixed strategies. Technical Report LSE-CDAM-2004-01, CDAM Research Report.

Yin, Z.; Korzhuk, D.; Kiekintveld, C.; Conitzer, V.; ; and Tambe, M. 2010. Stackelberg vs. Nash in security games: interchangeability, equivalence, and uniqueness. In *Proc. of The 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 1139–1146.