

Using Correlated Strategies for Computing Stackelberg Equilibria in Extensive-Form Games

Jiří Čermák¹, Branislav Bošanský^{1,2}, Karel Durkota¹, Viliam Lisý^{1,3}, Christopher Kiekintveld⁴

¹ Agent Technology Center, Faculty of Electrical Engineering, Czech Technical University in Prague

² Department of Computer Science, Aarhus University

³ Department of Computing Science, University of Alberta

⁴ Department of Computer Science, University of Texas at El Paso

jiri.cermak@agents.fel.cvut.cz, branislav.bosansky@agents.fel.cvut.cz, karel.durkota@agents.fel.cvut.cz,
lisy@ualberta.ca, cdkiekintveld@utep.edu

Abstract

Strong Stackelberg Equilibrium (SSE) is a fundamental solution concept in game theory in which one player commits to a strategy, while the other player observes this commitment and plays a best response. We present a new algorithm for computing SSE for two-player extensive-form general-sum games with imperfect information (EFGs) where computing SSE is an NP-hard problem. Our algorithm is based on a correlated version of SSE, known as Stackelberg Extensive-Form Correlated Equilibrium (SEFCE). Our contribution is therefore twofold: (1) we give the first linear program for computing SEFCE in EFGs without chance, (2) we repeatedly solve and modify this linear program in a systematic search until we arrive to SSE. Our new algorithm outperforms the best previous algorithms by several orders of magnitude.

Introduction

The roles of players in many games are often asymmetric. One example is the ability of one player (*the leader*) to commit to a strategy, to which the other players (*the followers*) react by playing their best response. In real-world scenarios, the leader can model a market leader with the power to set the price for items or services, or a defense agency committing to a security protocol to protect critical facilities. Optimal strategies for the players in such situations are described by the Strong Stackelberg Equilibrium (SSE) (Leitmann 1978; von Stengel and Zamir 2004). There are many examples of successful applications of SSE in infrastructure protection (Tambe 2011) as well as protecting the environment and wildlife (Fang, Stone, and Tambe 2015).

In most of the existing works, the game models are simplified and do not consider the sequential interaction among players (Pita et al. 2008; Tsai et al. 2009; Shieh et al. 2012; Jiang et al. 2013). One reason is computational complexity, since computing SSE is often NP-hard when sequential interaction is allowed (Letchford and Conitzer 2010). Another reason is the lack of practical algorithms. The only algorithm designed specifically for computing SSE in two-player imperfect-information general-sum extensive-form games

(EFGs) was only introduced recently (Bosansky and Cermak 2015) and formulates the problem as a mixed-integer variant of sequence-form linear program (referred to as BC15). However, the scalability of BC15 is limited as it contains a binary variable for each sequence of actions of the follower.

Our main contribution is a novel algorithm computing SSE in EFGs that offers a dramatic speed-up in computation time compared to BC15. The key idea behind our algorithm is in computing a variant of SSE, where the leader commits to correlated strategies – i.e., the leader can send signals to the follower, and the best response of the follower is to follow these signals. We use this variant to find the original SSE by systematically restricting which signals the leader can send to the follower. This variant of SSE has previously been studied in single step games (Conitzer and Korzhyk 2011), finite turn-based and concurrent-move games (Bosansky et al. 2015), infinite concurrent-move stochastic games (Letchford et al. 2012), and security games (Xu et al. 2015; Rabinovich et al. 2015; Durkota et al. 2015). Formally, it has been defined as Stackelberg Extensive-Form Correlated Equilibrium (SEFCE) in (Bosansky et al. 2015). While it was shown that the utility value for the leader in SSE cannot be closely approximated by SEFCE (Letchford et al. 2012), we show that one can use SEFCE to compute SSE. Since there was no previously known algorithm for computing SEFCE in EFGs, we also show that SEFCE can be found in polynomial time in EFGs without chance.

The paper is structured as follows. After introducing the formalism of EFGs, we formally define both SSE and SEFCE, and give an example of a game where these concepts differ. Next, we present the linear program (LP) for computing SEFCE in EFGs without chance (we describe a modified LP for EFGs with chance in the full version of the paper). Afterwards, we present our algorithm for computing SSE in EFGs that iteratively solves the LP for SEFCE with additional constraints until SSE is reached. Finally, we provide three variants of our algorithm and show that each variant significantly improves the computation time compared to BC15 on randomly generated games and an example of a search game.

Technical Background

Extensive-form games model sequential interactions between players and can be visually represented as game trees. Formally, a two-player EFG is defined as a tuple $G = (\mathcal{N}, \mathcal{H}, \mathcal{Z}, \mathcal{A}, u, \mathcal{C}, \mathcal{I})$: $\mathcal{N} = \{l, f\}$ is a set of players, the leader and the follower. We use i to refer to one of the players, and $-i$ to refer to his opponent. \mathcal{H} denotes a finite set of *nodes* in the game tree. Each node corresponds to a unique *history* of actions taken by all players and chance from the root of the game; hence, we use the terms history and node interchangeably. We say that h is a *prefix* of h' ($h \sqsubseteq h'$) if h lies on a path from the root of the game tree to h' . \mathcal{A} denotes the set of all actions. $\mathcal{Z} \subseteq \mathcal{H}$ is the set of all *terminal nodes* of the game. For each $z \in \mathcal{Z}$ we define a *utility function* for each player i ($u_i : \mathcal{Z} \rightarrow \mathbb{R}$). Chance player selects actions based on a fixed probability distribution known to all players. Function $\mathcal{C} : \mathcal{H} \rightarrow [0, 1]$ denotes the probability of reaching node h due to chance; $\mathcal{C}(h)$ is the product of chance probabilities of all actions in history h .

Imperfect observation of player i is modeled via *information sets* \mathcal{I}_i that form a partition over $h \in \mathcal{H}$ where i takes action. Player i cannot distinguish between nodes in any information set $I \in \mathcal{I}_i$. We overload the notation and use $A(I_i)$ to denote possible actions available in each node from information set I_i . We assume that action a uniquely identifies the information set where it is available. We assume *perfect recall*, which means that players remember history of their own actions and all information gained during the course of the game. As a consequence, all nodes in any information set I_i have the same history of actions for player i .

Pure strategies Π_i assign one action for each $I \in \mathcal{I}_i$. A more efficient representation in the form of *reduced pure strategies* Π_i^* assigns one action for each $I \in \mathcal{I}_i$ reachable while playing according to this strategy. A *mixed strategy* $\delta_i \in \Delta_i$ is a probability distribution over Π_i . For any pair of strategies $\delta \in \Delta = (\Delta_l, \Delta_f)$ we use $u_i(\delta) = u_i(\delta_i, \delta_{-i})$ for the expected outcome of the game for player i when players follow strategies δ . A *best response* of player i to the opponent's strategy δ_{-i} is a strategy $\delta_i^{BR} \in BR_i(\delta_{-i})$, where $u_i(\delta_i^{BR}, \delta_{-i}) \geq u_i(\delta'_i, \delta_{-i})$ for all $\delta'_i \in \Delta_i$.

Strategies in EFGs with perfect recall can be compactly represented by using the sequence form (Koller, Megiddo, and von Stengel 1996). A *sequence* $\sigma_i \in \Sigma_i$ is an ordered list of actions taken by a single player i in history h . \emptyset stands for the empty sequence (i.e., a sequence with no actions). A sequence $\sigma_i \in \Sigma_i$ can be extended by a single valid action a taken by player i , written as $\sigma_i a = \sigma'_i$. We say that σ_i is a *prefix* of σ'_i ($\sigma_i \sqsubseteq \sigma'_i$) if σ'_i is obtained by finite number (possibly zero) of extensions of σ_i . We use $\sigma_i(I_i)$ and $\sigma_i(h)$ to denote the sequence of i leading to I_i and h , respectively. We use the function $I_i(\sigma'_i)$ to obtain the information set in which the last action of the sequence σ'_i is taken. For an empty sequence, function $I_i(\emptyset)$ returns the information set of the root node. A mixed strategy of a player can now be represented as a *realization plan* ($r_i : \Sigma_i \rightarrow \mathbb{R}$). A realization plan for a sequence σ_i is the probability that player i will play σ_i under the assumption that the opponent plays to allow the actions specified in σ_i to be played. By

$g_i : \Sigma_l \times \Sigma_f \rightarrow \mathbb{R}$ we denote the *extended utility function*, $g_i(\sigma_l, \sigma_f) = \sum_{z \in \mathcal{Z} | \sigma_l(z) = \sigma_l \wedge \sigma_f(z) = \sigma_f} u_i(z) \mathcal{C}(z)$. If no leaf is reachable with pair of sequences σ , value of g_i is 0.

Solution Concepts in EFGs

Here we provide a formal definition of Strong Stackelberg Equilibrium (SSE) (e.g., in (Leitmann 1978)) and Stackelberg Extensive-Form Correlated Equilibrium (SEFCE) (Bosansky et al. 2015) and give the intuition on an example game.

Definition 1. A strategy profile $\delta = (\delta_l, \delta_f)$ is a Strong Stackelberg Equilibrium if δ_l is an optimal strategy of the leader given that the follower best-responds. Formally:

$$(\delta_l, \delta_f) = \arg \max_{\delta'_l \in \Delta_l, \delta'_f \in BR_f(\delta'_l)} u_l(\delta'_l, \delta'_f). \quad (1)$$

The SSE of the game in Figure 1 (the first utility in every leaf is for the leader, second for the follower) prescribes the leader to commit to playing g in h_4 , j in h_5 , and k in h_6 . The strategy of the follower is then to play a in h_1 and d in h_2 , leading to the expected utility of 1 for the leader.

In SEFCE we allow the leader to send signals to the follower and condition his strategy on sent signals. More specifically, the leader chooses $\pi_f^* \in \Pi_f^*$ as the recommendations for the follower according to SEFCE before the game starts. The actual recommendation to play some action $a \in A(I_f)$ is revealed to the follower only after he reaches I_f . Therefore, the follower only knows the past and current recommendations, and the probability distribution from which the recommendations are drawn in the future.

Definition 2. A probability distribution λ on reduced pure strategy profiles Π^* is called a Stackelberg Extensive-Form Correlated Equilibrium if it maximizes the leader's utility subject to the constraint that whenever play reaches an information set I where the follower can act, the follower is recommended an action a according to λ such that the follower cannot gain by unilaterally deviating from a in I and possibly in all succeeding information sets given the posterior on the probability distribution of the strategy of the leader, defined by the actions taken by the leader so far.

The middle table in Figure 1 represents the distribution λ forming the SEFCE of the example game (rows are labeled by Π_l^* , columns by Π_f^*). The leader chooses the signals to the follower to be either $\{a, c\}$ or $\{a, d\}$ based on the probability distribution depicted in the table. Afterwards, the corresponding column defines a valid mixed strategy for the leader and the signals the follower receives in his information sets. More specifically, the follower can receive either c or d in h_2 . When the follower receives the recommendation to play c , the leader commits to mix uniformly between g and h in h_4 and to play i in h_5 . When the follower receives d as the recommendation, the leader commits to playing g in h_4 and j in h_5 . Finally in h_1 the follower is recommended to play a , while the leader commits to play k in h_6 to ensure that the follower does not deviate from playing a . The expected utility of the leader is 1.5 in SEFCE. Note that SSE in this representation always corresponds to a table where only a single column of the follower has non-zero values.

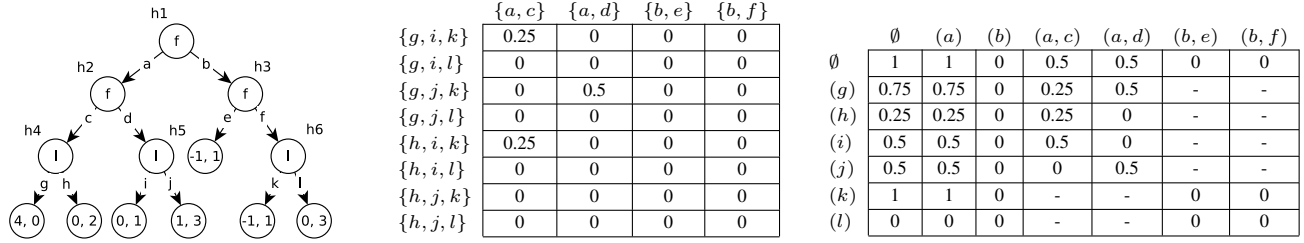


Figure 1: (Left) EFG with different SEFCE and SSE. (Middle) SEFCE distribution over \$\Pi^*\$. (Right) SEFCE correlation plan.

LP for Computing SEFCE

To compactly represent the behavior described in the middle table in Figure 1, we use a *correlation plan* (von Stengel and Forges 2008) that is quadratic in the size of the game tree, instead of the exponential representation using \$\Pi^*\$. A correlation plan for a sequence pair \$p(\sigma_l, \sigma_f)\$ represents the expected probability that \$\sigma_l\$ will be played if actions from \$\sigma_f\$ are recommended to the follower. In order to model SEFCE strategies using the correlation plan, the follower must be able to determine whether following the signal is the best response. Therefore, we need to specify the plan for so called *relevant sequences*. Consider our game in Figure 1; when the follower receives, for example, signal \$c\$ in \$h_2\$, the follower needs to know the commitment of the leader in the related part of the tree – i.e., in both nodes \$h_4, h_5\$ – to evaluate whether following the signal is the best response.

Definition 3. A pair of sequences \$(\sigma_l, \sigma_f)\$ is termed *relevant* if and only if either \$\sigma_f = \emptyset\$, or \$\exists h, h' \in \mathcal{H}, h' \sqsubseteq h; \sigma_l = \sigma_l(h) \wedge h' \in I_f(\sigma_f)\$.

By \$rel(\sigma_i)\$ we denote the set of sequences of \$-i\$ which form a relevant pair with \$\sigma_i\$. In our example \$rel((a)) = rel((b)) = \Sigma_l, rel((b, e)) = rel((b, f)) = \{\emptyset, (k), (l)\}\$, and \$rel((a, c)) = rel((a, d)) = \{\emptyset, (g), (h), (i), (j)\}\$.

Definition 4. A correlation plan (von Stengel and Forges 2008) is a partial function \$p : \Sigma_l \times \Sigma_f \to \mathbb{R}\$ such that there is a probability distribution \$\lambda\$ on the set of reduced strategy profiles \$\Pi^*\$ so that for each relevant sequence pair \$(\sigma_l, \sigma_f)\$, the term \$p(\sigma_l, \sigma_f)\$ is defined and fulfills \$p(\sigma_l, \sigma_f) = \sum_{(\pi_l, \pi_f) \in \Pi^*} \lambda(\pi_l, \pi_f)\$ where \$\pi_l, \pi_f\$ prescribe playing all of the actions in \$\sigma_l\$ and \$\sigma_f\$, respectively.

Let us now describe the SEFCE strategies (middle table in Figure 1) using the correlation plan (the right table; rows are labeled by \$\Sigma_l\$, columns by \$\Sigma_f\$). Every column of the table corresponds to an expected probability of the occurrence of the leader's sequences when the follower follows his recommendations in the future, and gets the recommendation to play the \$\sigma_f\$ corresponding to this column. We use '-' to mark irrelevant sequence pairs. The entry for every \$\sigma_l, \sigma_f\$ is the sum of all the entries corresponding to the pure strategies containing all the actions from \$\sigma_l\$ and \$\sigma_f\$ in the middle table of Figure 1. The behavior in columns corresponding to sequences \$(a, c)\$ and \$(a, d)\$ matches the behavior discussed in the previous section. The behavior in columns for sequence \$\emptyset\$ and \$(a)\$ corresponds to the expected probability of playing

sequences of the leader given the probability of recommendations, e.g., the probability of playing \$g\$ for the recommendation \$(a)\$ is equal to \$p((g), (a, c)) + p((g), (a, d))\$ (in this case the leader will play uniformly either \$g\$ with probability 0.5 according to the column for \$(a, c)\$ or \$g\$ with probability 1 according to the column for \$(a, d)\$). Probabilities in column for \$(a)\$ allow the follower to evaluate his choices in \$h_1\$.

Now we are ready to describe the LP for computing SEFCE in EFGs without chance that uses correlation plan:

$$\max_{p, v} \sum_{\sigma_l \in \Sigma_l} \sum_{\sigma_f \in \Sigma_f} p(\sigma_l, \sigma_f) g_l(\sigma_l, \sigma_f) \quad (2)$$

$$\text{s.t. } p(\emptyset, \emptyset) = 1; \quad 0 \leq p(\sigma_l, \sigma_f) \leq 1 \quad (3)$$

$$p(\sigma_l(I), \sigma_f) = \sum_{a \in A(I)} p(\sigma_l(I)a, \sigma_f) \quad \forall I \in \mathcal{I}_l, \forall \sigma_f \in rel(\sigma_l) \quad (4)$$

$$p(\sigma_l, \sigma_f(I)) = \sum_{a \in A(I)} p(\sigma_l, \sigma_f(I)a) \quad \forall I \in \mathcal{I}_f, \forall \sigma_l \in rel(\sigma_f) \quad (5)$$

$$v(\sigma_f) = \sum_{\sigma_l \in rel(\sigma_f)} p(\sigma_l, \sigma_f) g_f(\sigma_l, \sigma_f) + \sum_{I \in \mathcal{I}_f; \sigma_f(I) = \sigma_f} \sum_{a \in A_f(I)} v(\sigma_f a) \quad \forall \sigma_f \in \Sigma_f \quad (6)$$

$$v(I, \sigma_f) \geq \sum_{\sigma_l \in rel(\sigma_f)} p(\sigma_l, \sigma_f) g_f(\sigma_l, \sigma_f(I)a) + \sum_{I' \in \mathcal{I}_f; \sigma_f(I') = \sigma_f(I)a} v(I', \sigma_f)$$

$$\forall I \in \mathcal{I}_f, \forall \sigma_f \in \bigcup_{h \in I} rel(\sigma_l(h)), \forall a \in A(I) \quad (7)$$

$$v(\sigma_f(I)a) = v(I, \sigma_f(I)a) \quad \forall I \in \mathcal{I}_f, \forall a \in A(I) \quad (8)$$

The LP is derived from the computation of Extensive-Form Correlated Equilibrium (von Stengel and Forges 2008) by omitting the incentive constraints for the leader and maximizing the expected utility of the leader (this is similar to the approach in normal-form games (Conitzer and Korzhyk 2011)). Constraints (3) to (5) ensure that \$p\$ is a well-formed correlation plan. Constraint (6) ensures that \$v(\sigma_f)\$ represents the expected utility of playing \$\sigma_f\$ for the follower, when he follows his recommendations. The first sum represents the expected utility of the leaves reached by playing according to \$\sigma_l\$ and \$\sigma_f\$, the second sum represents the contribution of the expected utility from information sets reached by the continuations of \$\sigma_f\$. Constraint (7) ensures that \$v(I_f, \sigma_f)\$ is the maximum over all possible sequences leaving \$I_f\$ (denoted as \$\sigma_f(I_f)a\$ for all \$a \in A(I_f)\$) after the follower has received recommendation \$\sigma_f\$. Finally, constraint (8) forces the move recommended to the follower in \$I_f\$ to be optimal.

Definition 5. We say that p uses inconsistent recommendation in $I \in \mathcal{I}_f$ if and only if p defines two different recommendations for the follower in I . Formally, $\exists a, a' \in A(I), a \neq a', \exists \sigma_l, \sigma'_l \in \bigcup_{h \in I} \sigma_l(h) p(\sigma_l, \sigma_f(I)a) > 0 \wedge p(\sigma'_l, \sigma_f(I)a') > 0$. If there exists no such information set we say that p uses only consistent recommendations.

Theorem 1. Assume a solution of the LP as described in eqs. (2) to (8) such that there are only consistent recommendations for the follower. There is a SSE strategy that can be found in polynomial time from the p variables.

Proof. First, we show how the strategy of the leader is constructed. In every $I \in \mathcal{I}_l$ there is a subset Σ_r of relevant sequences $rel(\sigma_l(I))$ played with a positive probability. The behavior in I is specified by $p(\sigma_l(I)a, \sigma_f)$ for all $a \in A(I)$ for arbitrary $\sigma_f \in \Sigma_r$, as the behavior is the same for all $\sigma'_f \in \Sigma_r$. This is guaranteed by constraint (5) – it forces the probability of $p(\sigma_l, \sigma'_f)$ to be equal to the sum of $p(\sigma_l, \sigma''_f)$ over all extensions σ''_f of σ'_f . Since there can be only a single extension played with positive probability (assumed consistent recommendations) the probabilities must be equal.

For every $I \in \mathcal{I}_f$ there exists at most one action $a \in A(I)$ with $p(\sigma_l, \sigma_f a) > 0$ for some $\sigma_l \in \Sigma_l$ and $\sigma_f = \sigma_f(I)$ (consistent recommendations). By taking these actions and arbitrary actions in information sets where there is no such action a , we obtain a pure strategy for the follower. Finally, due to the correctness of the strategy of the leader proved in the previous step and constraints (6–8), this pure strategy is a best response of the follower. \square

Theorem 2. Assume a solution of the LP as described in eqs. (2) to (8). The objective value is greater than or equal to the expected utility of the leader in SSE.

Proof. Theorem 1 shows that in case the leader uses only consistent recommendations, the value of the LP corresponds to the expected utility of the leader in SSE. If the leader can also use inconsistent recommendations, the value of the LP can be only greater or equal. \square

Algorithm Computing SSE

In this section we describe the algorithm for computing SSE. The algorithm uses the linear program for computing SEFCE in case the game is without chance. Otherwise, a slightly modified version of this LP is required, however, the algorithm remains the same. Due to the space constraints, we describe this modified LP in details in the full version of the paper and refer to one of these two LPs as UB-SSE-LP.

The high level idea of our algorithm (depicted in Algorithm 1) is a recursive application of the following steps: (1) solve the UB-SSE-LP, (2) detect the set of information sets of the follower with inconsistent recommendations \mathcal{I}_{in} , (3) restrict the leader to use only consistent recommendations in \mathcal{I}_{in} by adding new constraints to the UB-SSE-LP. Restrictions are added cumulatively, until we arrive at a restricted UB-SSE-LP yielding only consistent p . The expected utility for the leader and the correlation plan p in this solution correspond to a candidate for the expected utility of the leader and the strategies of players in SSE. It is only

Input: An UB-SSE-LP P

Output: leader's expected utility and strategy profile in SSE

```

1  $M \leftarrow \{(\infty, \emptyset)\}; LB \leftarrow -\infty; p_c \leftarrow \emptyset$ 
2 while  $M \neq \emptyset$  do
3    $(UB, m) \leftarrow \max(M)$ 
4   if  $UB < LB$  then
5     | return  $(LB, p_c)$ 
6    $\text{apply}(m, P)$ 
7   if  $\text{feasible}(P)$  then
8     |  $(value, p) \leftarrow \text{solve}(P)$ 
9     |  $\mathcal{I}_{in} \leftarrow \text{inconsistentRecommendations}(p)$ 
10    | if  $\mathcal{I}_{in} = \emptyset$  then
11      | | if  $value > LB$  then  $LB \leftarrow value; p_c \leftarrow p$ 
12      | | else  $\text{addRestrictions}((UB, m), M, \mathcal{I}_{in}, value)$ 
13    |  $\text{revert}(m, P)$ 
14 return  $(LB, p_c)$ 

```

Algorithm 1: Algorithm for computing the SSE.

a solution candidate, since we have enforced actions, which may not be a part of SSE.

In more details, Algorithm 1 assumes as the input the UB-SSE-LP P for the game we want to solve. By *modification* $mod = (UB, m)$ we denote a pair of constraints m , to be added to P , and the upper bound UB on the value of P after adding the constraints. M is the set of modifications to be explored during the search for the SSE sorted in descending order of UB . The variable LB stores the expected utility for the leader in the best solution candidate found so far p_c . The main cycle of the algorithm starts on line 2, we iterate until there are no possible modifications of P left. On line 3 we remove the modification with the highest UB from M . We choose such $mod = (UB, m)$ in order to first explore modifications with the potential to lead to solution candidates with the highest expected utility for the leader. The algorithm verifies whether the modification mod (the one with the highest UB value) can improve the current best solution candidate (line 4). If not, the algorithm terminates and the best candidate found so far is the SSE. Otherwise, we add the constraints in m to P (line 6). If the modified P is feasible, the algorithm solves the LP (line 8) obtaining the expected utility of the leader and the correlation plan p . We find the set of information sets where the follower gets an inconsistent recommendation in p (line 9). If p uses only consistent recommendations, this solution corresponds to a solution candidate. If the expected utility for the leader is higher than for the best solution candidate found so far, we replace it (line 11). If \mathcal{I}_{in} is not empty, we generate new modifications to be applied to P and add them to M (line 12). The function addRestrictions will be discussed in more detail in the next subsection, as we explored several options of the modification generation. Finally, we revert the changes in m (line 13). If there are modifications left to be explored in M we continue with the next iteration. Every modification contains all the constraints added to the original P given as an input, therefore after revert we again obtain the original P .

When we enforce the whole strategy of the follower to be consistent with the SSE, the solution of the LP corresponds to the SSE, as it now maximizes the expected utility of the leader under the restriction the follower gets recom-

```

1 addRestrictions((UB, m), M,  $\mathcal{I}_{in}$ , value)
2    $I \leftarrow \text{getShallowest}(\mathcal{I}_{in})$ 
3   for  $a \in A(I)$  do
4      $UB_a \leftarrow \text{value}; m_a \leftarrow m$ 
5     for  $\sigma_l \in \text{rel}(\sigma_f(I)a)$  do
6        $m_a \leftarrow m_a \cup \{p(\sigma_l, \sigma_f(I)) = p(\sigma_l, \sigma_f(I)a)\}$ 
7        $M \leftarrow M \cup \{(UB_a, m_a)\}$ 

```

Algorithm 2: SI-LP.

```

1 addRestrictions((UB, m), M,  $\mathcal{I}_{in}$ , value)
2    $\mathcal{I}_c \leftarrow \text{chooseSets}(\mathcal{I}_{in})$ 
3    $UB' \leftarrow \infty; m' \leftarrow m$ 
4   for  $I \in \mathcal{I}_c$  do
5      $A_c \leftarrow \{a \in A(I) \mid \exists \sigma_l \in \Sigma_l p(\sigma_l, \sigma_f(I)a) > 0\}$ 
6     for  $a \in A_c$  do
7       for  $\sigma_l \in \text{rel}(\sigma_f(I)a)$  do
8          $m' \leftarrow m' \cup \{p(\sigma_l, \sigma_f(I)a) \leq b_a\}$ 
9          $m' \leftarrow m' \cup \{\sum_{a \in A_c} b_a = 1\}$ 
10     $M \leftarrow M \cup \{(UB', m')\}$ 

```

Algorithm 3: MILP.

mended the pure best response to the strategy of the leader. It remains to be shown, that we are guaranteed to add modifications to M which force the correct set of actions of the follower in every version of *addRestrictions*. The final correctness arguments will be provided after discussing the *addRestrictions* versions used.

Rules for Restricting Follower’s Behavior

We examine different approaches for method *addRestrictions* that generates new modifications of UB-SSE-LP resulting in three different variants of our new algorithm.

In Algorithm 2 we describe the first version of the function *addRestrictions*, labeled SI-LP. On line 2 we find the shallowest information set I , where the follower receives an inconsistent recommendation. We generate modification for every $a \in A(I)$. Every such modification enforces corresponding $a \in A(I)$ to be recommended deterministically. The upper bound is set to the value of P computed before invoking *addRestrictions*. The shallowest information set is chosen to avoid unnecessary modifications in deeper parts of the game tree, which might end up not being visited at all. This version of *addRestrictions* transforms Algorithm 1 to branch and bound algorithm.

By adding $mod = (UB, m)$ to M for every action in the shallowest information set with inconsistent recommendations, until no such set exists, we ensure that all of the actions which might form a part of SSE will be tried. The behavior in information sets with consistent recommendation need not be restricted, as we are sure that the follower has no incentive to deviate and therefore plays his best response maximizing the expected utility of the leader. Finally, since we assign to UB a value which forms an upper bound on the solution of P after adding constraints m , we are sure that if we terminate the algorithm on line 4 in Algorithm 1, there is indeed no possibility to encounter a solution candidate better than the one yielding the current LB.

A second option, presented in Algorithm 3, chooses $\mathcal{I}_c \subseteq$

\mathcal{I}_{in} (line 2) and restricts the recommendations in every $I \in \mathcal{I}_c$. The restriction in I is done in the following way. First, detect the subset A_c of $A(I)$ of actions which are recommended with positive probability (line 5) and make the recommendation mutually exclusive using binary variables (lines 8 and 9), converting the LP P to a mixed integer linear program (MILP). We use two options of creating \mathcal{I}_c . First, we create a singleton containing only the shallowest $I \in \mathcal{I}_{in}$, we refer to this algorithm as SI-MILP. Second, we let $\mathcal{I}_c = \mathcal{I}_{in}$, we refer to this algorithm as AI-MILP. Algorithm 1 using both SI-MILP and AI-MILP closely resembles constraint generation, with the difference that additional binary variables are also added in every iteration.

If we introduce a binary variable for every action of the follower in the game, we are guaranteed to obtain the SSE, as the MILP then finds a strategy profile maximizing the expected utility of the leader (ensured by the objective), while the follower plays a pure best response to the strategy of the leader (breaking ties in favor of the leader due to the objective), which is the definition of the SSE. If we create some partial enforcement of consistent recommendations using the binary variables and we obtain a pure strategy for the follower then this is again SSE, since the enforcement in the rest of the game would not make any difference as the follower already gets consistent recommendations there. Finally, since we restrict the follower’s recommendations until consistent recommendations are obtained, both MILP based rules indeed guarantee to find the SSE.

Experimental Evaluation

We now turn to the experimental evaluation of the three described variants of our algorithm for computing an SSE. We use BC15 (Bosansky and Cermak 2015) as a baseline, state-of-the-art algorithm for computing SSE in EFGs. Single-threaded IBM CPLEX 12.5 solver was used to compute all the (MI)LPs. We use two different domains previously used for evaluation of BC15: a search game representing the scenario where security units defend several targets against an attacker, and randomly generated games.

Search Game. The search game is played on a directed graph (see Figure 2). The follower aims to reach one of the destination nodes (D1 – D3) from starting node (E) in a given number of moves, while the leader aims to catch the follower with one of the two units operating in the shaded areas of the graph (P1 and P2). The follower receives different reward for reaching different destination node (the reward is randomly selected from the interval $[1, 2]$). The leader receives positive reward 1 for capturing the follower. Once the follower runs out of moves without reaching any goal or being captured, both players receive 0. The follower leaves tracks in the visited nodes that can be discovered if the leader visits the node. The follower can erase the tracks in the current node (it takes one turn of the game). The follower does not know the position of the patrolling units, the leader observes only the tracks left by the follower.

Randomly Generated Games. We use randomly generated games, where in each state of the game the number of

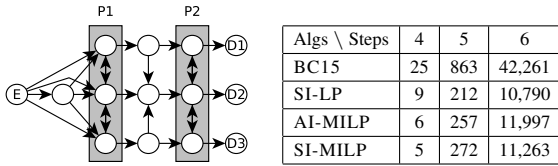


Figure 2: (Left) Search game graph. (Right) Runtimes in seconds for the search game with increasing depth.

Table 1: Number of games solved in given time intervals.

Algs \ Runtime	1s	5s	30s	2min	25min	4h
BC15	2	12	137	245	393	55
SI-LP	583	191	54	12	4	0
AI-MILP	529	259	51	5	0	0
SI-MILP	483	279	72	9	1	0

available actions is randomly generated up to a given parameter $\{2, \dots, \max_A\}$. Each action leads to a state where the opponent is to move and also generates an *observation* for the opponent. An observation is a number from a set $\{1, \dots, \max_O\}$ and determines partitioning of the nodes into the information sets – for player i , the nodes h with the same history of moves $\sigma_i(h)$ and the observations generated by the actions of the opponent $-i$ belong to the same information set. We generate games of differing sizes by varying parameters $\max_A = \{3, 4\}$, $\max_O = \{2, 3\}$, and depth of the game (up to 5 actions for each player). The utility for the players is randomly generated in the interval $[-100, 100]$. The utilities are correlated with factor set to -0.5 (1 represents identical utilities, -1 zero-sum utilities).

Results

The runtime results on random games are depicted in the top graph of Figure 3. The x-axis shows the number of realization plans of the follower, while the y-axis depicts the time in seconds needed to solve a given instance (both axes are logarithmic). The number of realization plans of the follower is a good estimate of how difficult the game is to solve as it reflects both the size of the game as well as the structure of information sets. Each point represents the mean time needed for every algorithm to solve the instances from a given time interval (at least 600 different instances). The standard errors of the mean values are very small compared to the differences between algorithms and not visible in the graph.

The results show that each of the variants significantly outperforms the previous state-of-the-art algorithm BC15. It typically takes around 10 minutes for BC15 to solve games with 10^6 realization plans of the follower, while our algorithms were often able to find solutions under a second. AI-MILP performs best on average, since it is the least sensitive to the different structure of the instances solved. AI-MILP fixes the behavior in a higher number of information sets and so it performs a successful trade off between the complexity of solving a single MILP and the number of MILP invocations. The second best approach on average is the SI-MILP. The average performance is slightly worse due to a higher number of MILP invocations needed to solve more difficult

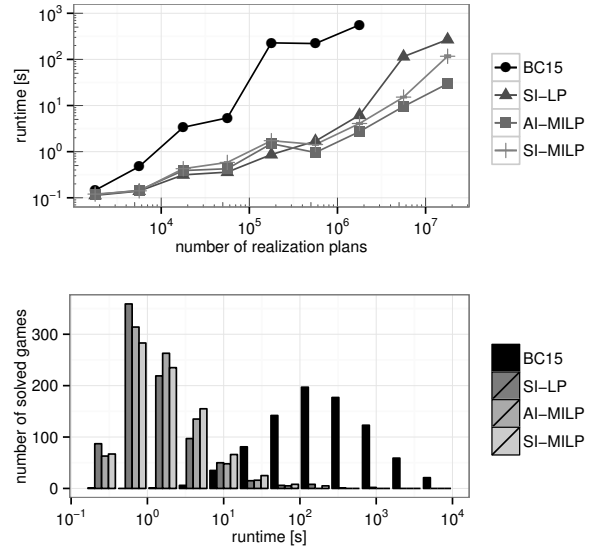


Figure 3: (Top) Runtimes on randomly generated games. (Bottom) Number of solved games in given time intervals.

instances (i.e., the ones with many inconsistent recommendations). Finally, the SI-LP has the worst average performance out of the variants of our new algorithm, since SI-LP needs even more LP invocations on more difficult instances in comparison to the previous variants of our algorithm.

Additionally, we provide in the bottom graph of Figure 3 a histogram for number of instances solved (y-axis) within a time interval (x-axis). The results were calculated on random games with the number of realization plans from interval $[3 \cdot 10^5, 3 \cdot 10^6]$. Despite a slightly worse average performance of SI-LP on these instances, it solved the highest number of instances very fast compared to the other two variants (SI-LP solves 69% of the instances under 1 second, while AI-MILP solves 62% and SI-MILP 57%). The histogram also shows the reason behind the worse average performance of SI-LP. There are multiple instances that SI-LP solves in more than 200 seconds, while such outliers are not present for the latter two variants (the worst outlier for SI-LP took 773 seconds, while the longest time for SI-MILP was 146 seconds, for AI-MILP 57 seconds and for BC15 2.5 hours). For clarity we provide the same data in coarser intervals in Table 1, where the outliers are clearly visible (the label of column represents the upper bound of the corresponding time interval, the label of the column to the left the lower bound of the time interval). The results show that SI-LP is more efficient on instances where the advantage in using the correlated strategies is marginal and there are only few information sets with inconsistent recommendations. On the other hand, AI-MILP offers higher robustness across different instances.

Finally, in Figure 2 we present the results on the search game. All the new approaches performed similarly, outperforming the BC15 in every setting. This shows that our algorithm outperforms BC15 even in an unfavorable setting. The search game has a specific structure, where the strategy

space of the leader is large (joint actions of the patrolling units), while the strategy space of the follower is significantly smaller. This structure is favorable for BC15, since it implies a relatively small number of binary variables (the MILP contains one binary variables for each sequence of the follower), with the overall size of the MILP being linear in the size of the game, while the size of our underlying LP is quadratic due to the correlation plan.

Conclusion

We present a novel domain-independent algorithm for computing Strong Stackelberg Equilibria (SSE) in extensive-form games that uses the correlated variant of Stackelberg Equilibria (Stackelberg Extensive-Form Correlated Equilibrium). This work opens several areas for future research. First, our algorithm can be adapted and applied for solving specific domains since its scalability is significantly better in comparison to the existing algorithms. Second, the scalability can most-likely be further improved by employing iterative approaches for solving the underlying linear program. Third, several questions were not addressed in our approach and remain open: Is it possible to generalize presented algorithm for computing SSE with multiple followers? Can we relax the assumption of perfect recall?

Acknowledgments

This research was supported by the Czech Science Foundation (grant no. 15-23235S), by the Danish National Research Foundation and The National Science Foundation of China (under the grant 61361136003) for the Sino-Danish Center for the Theory of Interactive Computation and Office of Naval Research Global (grant no. N62909-13-1-N256). This material is based upon work supported by the National Science Foundation (grant no. IIS-1253950). This research was supported by Alberta Innovates Technology Futures through the Alberta Innovates Centre for Machine Learning and Reinforcement Learning and AI Lab and by the Grant Agency of the Czech Technical University in Prague, grant No. SGS15/205/OHK3/3T/13 and SGS15/206/OHK3/3T/13.

References

- Bosansky, B., and Cermak, J. 2015. Sequence-Form Algorithm for Computing Stackelberg Equilibria in Extensive-Form Games. In *AAAI Conference on Artificial Intelligence*, 805–811.
- Bosansky, B.; Branzei, S.; Hansen, K. A.; Miltersen, P. B.; and Sorensen, T. B. 2015. Computation of Stackelberg Equilibria of Finite Sequential Games. In *11th Conference on Web and Informatics (WINE)*.
- Conitzer, V., and Korzhyk, D. 2011. Commitment to Correlated Strategies. In *AAAI Conference on Artificial Intelligence*.
- Durkota, K.; Lisy, V.; Bosansky, B.; and Kiekintveld, C. 2015. Approximate solutions for attack graph games with imperfect information. In *Decision and Game Theory for Security*, 228–249. Springer.
- Fang, F.; Stone, P.; and Tambe, M. 2015. When Security Games Go Green: Designing Defender Strategies to Prevent Poaching and Illegal Fishing. In *Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI*, 2589–2595.
- Jiang, A. X.; Yin, Z.; Zhang, C.; Tambe, M.; and Kraus, S. 2013. Game-theoretic Randomization for Security Patrolling with Dynamic Execution Uncertainty. In *12th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 207–214.
- Koller, D.; Megiddo, N.; and von Stengel, B. 1996. Efficient Computation of Equilibria for Extensive two-person Games. *Games and Economic Behavior* 247–259.
- Leitmann, G. 1978. On generalized Stackelberg strategies. *Journal of Optimization Theory and Applications* 26(4):637–643.
- Letchford, J., and Conitzer, V. 2010. Computing Optimal Strategies to Commit to in Extensive-Form Games. In *11th ACM conference on Electronic commerce*, 83–92.
- Letchford, J.; MacDermed, L.; Conitzer, V.; Parr, R.; and Isbell, C. L. 2012. Computing Optimal Strategies to Commit to in Stochastic Games. In *AAAI Conference on Artificial Intelligence*.
- Pita, J.; Jain, M.; Marecki, J.; Ordóñez, F.; Portway, C.; Tambe, M.; Western, C.; Paruchuri, P.; and Kraus, S. 2008. Deployed ARMOR protection: the application of a game theoretic model for security at the Los Angeles International Airport. In *7th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 125–132.
- Rabinovich, Z.; Jiang, A. X.; Jain, M.; and Xu, H. 2015. Information Disclosure as a Means to Security. In *14th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 645–653.
- Shieh, E.; An, B.; Yang, R.; Tambe, M.; Baldwin, C.; DiRenzo, J.; Maule, B.; and Meyer, G. 2012. Protect: A deployed game theoretic system to protect the ports of the united states. In *11th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 13–20.
- Tambe, M. 2011. *Security and Game Theory: Algorithms, Deployed Systems, Lessons Learned*. Cambridge University Press.
- Tsai, J.; Kiekintveld, C.; Ordóñez, F.; Tambe, M.; and Rathi, S. 2009. IRIS - A Tool for Strategic Security Allocation in Transportation Networks Categories and Subject Descriptors. In *8th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 37–44.
- von Stengel, B., and Forges, F. 2008. Extensive-form Correlated Equilibrium: Definition and Computational Complexity. *Mathematics of Operations Research* 33(4):1002–1022.
- von Stengel, B., and Zamir, S. 2004. Leadership with Commitment to Mixed Strategies. Technical report, CDAM Research Report LSE-CDAM-2004-01.
- Xu, H.; Rabinovich, Z.; Dughmi, S.; and Tambe, M. 2015. Exploring Information Asymmetry in Two-Stage Security Games. In *AAAI Conference on Artificial Intelligence*.

Appendix

Extensive-Form Games with Chance

In this section we analyze the situation when EFGs contain chance. With chance, one cannot directly use the correlation plan as defined in Definition 4 to compute SEFCE – von Stengel and Forges showed (2008) that the presence of chance disrupts the structure of relevant sequences. More specifically, by using correlation plan p , one can describe a richer set of strategies that does not correspond to Extensive-Form Correlated Equilibrium, defined as a probability distribution over (reduced) pure strategy profiles.

We are, however, interested in finding SSE for a given EFG with chance – i.e., such a solution of the LP in eqs. (2) to (8) which recommends only a single pure strategy to the follower. We can exploit this restriction and modify the LP that uses correlation plan according to Definition 4. We show that this new LP (1) returns SSE if all recommendations for the follower are consistent; (2) computes an upper bound on the expected utility of the leader in SSE.

The main difference caused by chance is that multiple nodes in information sets of the leader can now be reached with a non-zero probability in p even if the recommendation for the follower is always consistent. This is not the case in games without chance, since each node in an information set of the leader has the same history for the leader (due to *perfect recall*) and a different sequence of actions of the follower. Therefore, we need to add constraints ensuring that the strategy of the leader cannot depend on actions of chance unobserved by the leader (i.e., they spawn different nodes in an information set of the leader).

Definition 6. Let $I \in \mathcal{I}_l$ be an information set of the leader and $h \in I$ nodes in this information set. Define $\Gamma_I = (\gamma_1, \gamma_2, \dots)$ to be a partitioning of the nodes in I based on the history of chance. Every node h belongs to exactly one set γ_k such that all nodes in γ_k share the same sequence of chance actions. Formally,

$$\forall h \in I \exists \gamma_k \in \Gamma_I h \in \gamma_k$$

$$\forall h, h' \in \gamma_k \sigma_C(h) = \sigma_C(h')$$

$$\forall h \in \gamma_k \forall h' \in \gamma_j k \neq j \wedge \sigma_C(h) \neq \sigma_C(h')$$

Next, we extend the definition of relevant sequences to a restricted set such that the sequences are relevant due to a specific node in the information set of the leader h :

$$rel(h) = \{\sigma_f \mid \exists h' \in \mathcal{H}, h' \sqsubseteq h; h' \in I_f(\sigma_f)\},$$

additionally there can be no pair of sequences $\sigma_f, \sigma'_f \in rel(h)$ such that one is a strict prefix of another: $\sigma_f \sqsubseteq \sigma'_f$ and $\sigma_f \neq \sigma'_f$. Now we can give a modified version of LP used in our algorithm:

$$\max_{p, v} \sum_{\sigma_l \in \Sigma_l} \sum_{\sigma_f \in \Sigma_f} p_{\sigma_l, \sigma_f} g_l(\sigma_l, \sigma_f) \quad (9)$$

$$\text{s.t. constraints (3)-(8)} \quad (10)$$

$$s_{I, a_l} = \sum_{h \in \gamma_k} \sum_{\sigma_f \in rel(h)} p(\sigma_l(I) a_l, \sigma_f) \quad \forall I \in \mathcal{I}_l \forall a_l \in A_l(I) \forall \gamma_k \in \Gamma_I \quad (11)$$

$$0 \leq s_{I, a_l} \leq 1 \quad \forall I \in \mathcal{I}_l \forall a_l \in A_l(I) \quad (12)$$

We added a new variable s_{I, a_l} that represents the realization probability of action a_l being played in the information set I . Constraint (11) now ensures that the strategy in this information set cannot depend on chance, as shown by the following Lemma and Theorem.

Lemma 1. Assume a solution of the LP described in eqs. (9-12) such that there are only consistent recommendations for the follower. Then, for every information set of the leader $I \in \mathcal{I}_l$ and every partition $\gamma_k \in \Gamma_I$ there exists at most one sequence of the follower $\sigma_f \in rel(h)$ for some node $h \in \gamma_k$ such that $p(\sigma_l(h), \sigma_f) > 0$.

Proof. For contradiction, let σ_f, σ'_f be such different relevant sequences for which $p(\sigma_l(I), \sigma_f) > 0$ and $p(\sigma_l(I), \sigma'_f) > 0$ and let $h, h' \in \gamma_k$ (not necessarily different), such that $\sigma_f \in rel(h)$ and $\sigma'_f \in rel(h')$.

Now, thanks to constraint (5) we know that there must be a predecessor of these nodes with strictly positive value of p . Let σ''_f be the longest common prefix of the two sequences of the follower σ_f and σ'_f . By construction of $rel(h)$ we now that $\sigma_f \neq \sigma'_f \neq \sigma''_f$. We distinguish two cases: (1) either the first different action in sequences σ_f and σ'_f after the longest common prefix σ''_f is taken in the same information set, or (2) there are two different information sets in which these different actions are taken. The first case contradicts the assumption that there are no inconsistent recommendations for the follower. In the second case, these two information sets must be reached either due to a different action of the leader or chance, since σ''_f is the longest common prefix and the first different actions in the sequences of the follower are played in these information sets. However, this contradicts the fact that the history for the leader is the same due to the assumption of perfect recall, as well as the history of chance is the same for these two nodes due to the partitioning of Γ_I . \square

Theorem 3. Assume a solution of the LP as described in eqs. (9-12) such that there are only consistent recommendations for the follower. There is a SSE strategy that can be found in polynomial time from variables p .

Proof. We need to show how $p(\sigma_l, \sigma_f)$ is translated to a valid mixed strategy that corresponds to SSE. First, we show how the strategy of the leader is built. There is only a single relevant sequence $\sigma_f \in \bigcup_{h \in \gamma_k} rel(h)$ with a strictly positive probability in every information set of the leader I (thanks to Lemma 1). Since the history of actions of the leader is the same for all nodes in this information set (perfect recall) and thanks to the constraint (5), the value of the positive probability is the same for each relevant sequence of the follower (only the prefixes of σ_f have positive probability) for each partition γ_k . Therefore, constraint (11) now ensures that the behavior of the leader is unique for all relevant sequences of the follower. And so, for every information set of the leader variables p define a valid strategy that maximizes the expected utility of the leader and that can be computed in the same way as in the standard sequence form in EFGs (e.g., (Koller, Megiddo, and von Stengel 1996)).

For the case of the follower, there is no inconsistent recommendations by assumption; hence, for every information set I of the follower there exists at most one action $a_f \in A(I)$ with $p(\sigma_l, \sigma_f a_f) > 0$ for some $\sigma_l \in \Sigma_l$ and $\sigma_f = \sigma_f(I)$. By taking these actions and arbitrary actions in information sets where there is no such action a_f , we obtain a pure strategy for the follower. Finally, thanks to the constraints (6–8), this pure strategy is a best response of the follower to the strategy of the leader obtained in the previous step. \square

Theorem 4. *Assume a solution of the LP as described in eqs. (9-12). The objective value of the LP is greater or equal to the expected utility of the leader in SSE.*

Proof. The theorem holds as a direct consequence of Theorem 3. We know that in case there are no inconsistent recommendations for the follower, the value of the objective corresponds to SSE. If we do not restrict the recommendations, the value can be only greater or equal. \square