



HAL
open science

Stationary Strong Stackelberg Equilibrium in Discounted Stochastic Games

Víctor Bucarey, Eugenio Della Vecchia, Alain Jean-Marie, Fernando Ordóñez

► **To cite this version:**

Víctor Bucarey, Eugenio Della Vecchia, Alain Jean-Marie, Fernando Ordóñez. Stationary Strong Stackelberg Equilibrium in Discounted Stochastic Games. [Research Report] RR-9271, Inria. 2021, pp.62. hal-02144095v3

HAL Id: hal-02144095

<https://inria.hal.science/hal-02144095v3>

Submitted on 12 Mar 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Stationary Strong Stackelberg Equilibrium in Discounted Stochastic Games

Víctor Bucarey, Eugenio Della Vecchia, Alain Jean-Marie, Fernando
Ordóñez

**RESEARCH
REPORT**

N° 9271

May 2019

Project-Teams NEO and INOCS



Stationary Strong Stackelberg Equilibrium in Discounted Stochastic Games

Víctor Bucarey^{*†}, Eugenio Della Vecchia[‡], Alain Jean-Marie[§],
Fernando Ordóñez[¶]

Project-Teams NEO and INOCS

Research Report n° 9271 — version 3 — initial version May 2019 —
revised version March 2021 — 62 pages

Abstract: In this work we focus on Stackelberg equilibria for discounted stochastic games. We begin by formalizing the concept of Stationary Strong Stackelberg Equilibrium (SSSE) policies for such games. We provide classes of games where the SSSE exists, and we prove via counterexamples that SSSE does not exist in the general case. We define suitable dynamic programming operators whose fixed points are referred to as Fixed Point Equilibrium (FPE). We show that the FPE and SSSE coincide for a class of games with Myopic Follower Strategy. We provide numerical examples that shed light on the relationship between SSSE and FPE and the behavior of Value Iteration, Policy Iteration and Mathematical programming formulations for this problem. Finally, we present a security application to illustrate the solution concepts and the efficiency of the algorithms studied in this report.

Key-words: Stochastic games, Stackelberg Equilibrium, Optimal control

Version note: This third version of the document: fixes minor issues of form and typos; adds statements to Theorem 1 and Lemma 3.

* Département d'Informatique, Université Libre de Bruxelles, Brussels, Belgium.

† Inria Lille-Nord Europe, Villeneuve d'Ascq, France. vbucarey@ulb.ac.be

‡ Departamento de Matemática, Universidad Nacional de Rosario, Argentina. eugenio@fceia.unr.edu.ar

§ Inria, LIRMM, University of Montpellier, CNRS, Montpellier, France. alain.jean-marie@inria.fr

¶ Departamento de Ingeniería Industrial, Universidad de Chile, Chile. fordon@dii.uchile.cl

**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Équilibre de Stackelberg Stationnaire Fort dans les jeux stochastiques actualisés

Résumé : Dans cet article, nous nous focalisons sur les équilibres de Stackelberg pour les jeux stochastiques actualisés. Nous commençons par formaliser le concept d'équilibre fort de Stackelberg en politiques stationnaires (*Strong Stationary Stackelberg Equilibria*, SSSE) pour ces jeux. Nous exhibons des classes de jeux pour lesquels le SSSE existe et nous montrons par des contre-exemples que les SSSE n'existent pas dans le cas général. Nous définissons des opérateurs de programmation dynamique appropriés pour ce concept dont les points fixes sont nommés FPE (*Fixed Point Equilibria*). Nous montrons que le FPE et le SSSE coïncident pour la classe de jeux stochastiques avec stratégie du "follower" myope (*Myopic Follower Strategy*, MFS). Nous montrons des exemples numériques qui éclairent la relation entre SSSE et FPE ainsi que le comportement des algorithmes Value Iteration, Policy Iteration, et la formulation par Programmation Mathématique de ce problème. Finalement, nous décrivons une application dans le domaine de la sécurité pour illustrer les concepts de solution et l'efficacité des algorithmes introduits dans ce rapport.

Mots-clés : Jeux stochastiques, Équilibre de Stackelberg, Contrôle optimal

1 Introduction

Stackelberg games model interactions between strategic agents, where one agent, the leader, can enforce a commitment to a strategy and the remaining agents, referred to as followers, take that decision into account when selecting their own strategies. This Stackelberg game interaction can be extended to a multistage setting where leader and followers repeatedly make strategic decisions. Such dynamic Stackelberg models have been considered in applications in economics [1], marketing and supply chain management [12], dynamic congestion pricing [19], and security [5]. For example, in a dynamic security application, a defender could decide on a strategy to patrol a number of targets over multiple periods and the attackers would take this defender patrol into consideration when deciding whether to attack and where in each period.

The sequential interaction between leader and follower in Stackelberg games and Stackelberg stochastic games, characterizes its equilibrium solutions with a system of optimality conditions that are in general non-linear and non convex. Previous work has developed specialized methods to find such equilibrium solutions for different problem types. In the specific case of Stackelberg security games, bi-level mathematical optimization formulations have been used to compute the Strong Stackelberg Equilibrium solution [9] and, in the case of stochastic games, the Strong Stackelberg Equilibrium in *stationary strategies* (SSSE) [5]. An alternative method determines the SSSE by solving a non-linear potential game formulation [3]. These solution methods are either tailored for specific problems or only capable of solving small instances. It is therefore important to study alternative solution methods for Stackelberg stochastic games that are general and could solve large instances.

In this paper, we propose to use *iterative* algorithms, based on the operator formalism, to compute the SSSE solution of stochastic games. Our proposal is to use the well-known algorithms for solving Markov Decision Processes, Value Iteration and Policy Iteration, based on a suitably defined *dynamic programming operator*. Such algorithms are both conceptually simple and have a small computational burden per iteration. However, their use raises the question of convergence: do they converge, and if so, do they converge to some SSSE? Answering these questions led us to: a) conclude that SSSE do not necessarily exist, something not obvious from the current literature; b) identify classes of games where iterative algorithms converge to an SSSE. We then exploited this property in the analysis of a model of dynamic planning of police patrols on a transportation graph.

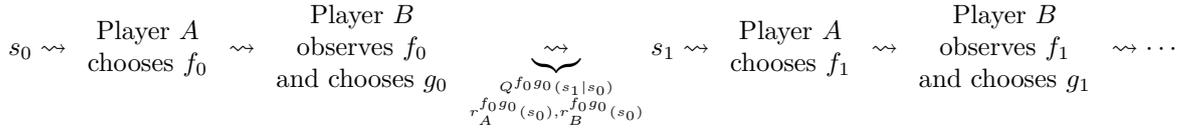
In the remainder of this introduction we clarify the problem under consideration, present related literature, describe the contributions of our work and introduce the notation that is used in this paper.

1.1 Problem Statement

We present now the notion of Strong Stackelberg Equilibria in Stationary policies (SSSE) for Stochastic Games.

Consider a dynamic system evolving in discrete time on a finite set of states, where

two players control the evolution. Players have a perfect information on the state of the system. One of them, called *Leader* or Player A, observes the current state s and commits to a, possibly mixed, strategy f depending solely on the state s . Then the other player, called *Follower* or Player B, observes the state and *strategy* of Player A and plays his best response denoted by g . Given the selected strategies there is a one-step reward for each player ($r_A^{fg}(s)$ and $r_B^{fg}(s)$ for player A and B, respectively) and a random transition probability $Q^{fg}(s'|s)$ to another state s' . This dynamics is illustrated in the following scheme:



Aggregated payoffs for both players are evaluated with the expected total discounted revenue over the infinite horizon, each player having their own discount factor. The aim for the leader is to find a policy that, in each state, maximizes her revenue taking into account that the follower will observe this policy and will respond by optimizing his own payoff. This general “Stackelberg” approach to the solution of the game is complemented with the rule that if the follower is indifferent between strategies, he chooses the one that benefits the leader: this refinement is the *strong* Stackelberg solution.

1.2 Related bibliography

The study of Strong Stackelberg Equilibria (SSE) has received much attention in the recent literature due to its relevance in security applications [9]. In static games, the need to generalize the standard Stackelberg equilibrium has been pointed out by Leitman [10] who introduced a conservative version of it. This generalization is formalized in Breton *et al* [2] as a Weak Stackelberg Equilibrium, together with the definition of the optimistic generalization, the Strong Stackelberg Equilibrium. The relationship between SSE and bi-level optimization appears in [2], solution methods for static Stackelberg Games are discussed in [6].

Stackelberg equilibria in multi-stage and dynamic games have been studied in [14, 15]. In particular, authors propose in [14] to focus on *feedback* strategies that can be obtained via dynamic programming. The idea is reused in [2] which introduces Strong *Sequential* Stackelberg Equilibria, in a setting very similar to ours. The notable difference is that, in the problem they consider, the follower gets to observe the *action* of the leader, not just its strategy. In their analysis, the formalism of operators linked to dynamic programming, introduced by Denardo [7] and developed in [18], is essential. Our analysis uses this formalism as well.

The stochastic game model we study in this paper is also the topic of Vorobeychik and Singh in [16]. Although they do not provide any formal definition of SSSE, authors show that SSSE always exist in stochastic games for *team* games where both players have the same rewards by reducing the game to a MDP. They also propose mathematical

programming formulations to find the SSSE, extending the analysis for MDPs (see [13, ch. 6.]) and Nash equilibrium in stochastic games [8], for this case. Similar mathematical programming formulations are established in [5] and [17] for problems in security applications. However, no prior work has provided a proof of the relationship between the solutions of these mathematical programming formulations and the SSSE of the stochastic game being considered. In this paper we present conditions that guarantee the existence of SSSE for stochastic games in diverse classes of problems, including team games. Furthermore we present numerical examples that suggest that the mathematical programming formulation computes the SSSE solution when it exists.

The complexity of computing a SSE is studied in [11]. That work shows, by reduction to 3SAT, that it is NP-hard to determine a SSE for a Stackelberg Stochastic Game with any discount factor $\beta > 0$ common for both players. The possibility that a Stackelberg Stochastic Game does not have a SSE is not mentioned.

As mentioned above, security applications are an important motivation for research on dynamic games. An attacker-defender Stackelberg security game is also considered in [3] for a repeated stochastic Markov chain game. This problem is represented as a potential game in terms of a suitable Lyapunov function, which is used to prove convergence results to compute the strong Stackelberg equilibrium [4]. While this method is a general approach for Ergodic Markov chains, the results presented show solutions only for instances with few states. It is interesting to single out the stochastic game model described in [8, Chapter 6.3]. The authors are interested in the average reward and the solution concept used is the Nash Equilibrium, a choice different from ours. However, the model has the feature that only one player, the defender, controls the transitions between states. This feature is one of the properties that guarantees the existence of SSSE, as we will show later.

1.3 Contribution

While previous work has formulated Stackelberg equilibrium for stochastic games and considered different solution methods to compute the SSSE, to the best of our knowledge the general question of the existence of SSSE is still largely open, in the sense that, so far: a) no case of non-existence is reported, b) few sufficient conditions for existence have been established. Furthermore, no work to date has advocated the use of the operator method, neither for the mathematical analysis of the problem, nor for its algorithmic solution.

We contribute to the issue in the following ways: First, we give a formal definition of the Strong Stationary Stackelberg Equilibrium (SSSE) in stochastic games (Section 1.4). We develop the operator-based analysis of such games by introducing an operator acting on the space of value functions. The operator introduced is related to the one-step evaluation of each player's payoff. We then define *Fixed Point Equilibria* (FPE) as the fixed points of this operator. Next, we introduce the class of games with *Myopic Follower Strategy* (MFS), for which specific operators are relevant. We prove that these operators are contractive. Finally, we introduce the algorithms for computing FPE, for general games and for games with MFS. We prove the convergence of both Value Iteration

and Policy Iteration to the FPE of games with MFS. We also recall the Mathematical Programming formulation for SSSE. This is the topic of Section 2.

Next, we focus on the general question of existence of SSSE and FPE, and how they are related. We show that games with MFS and Team Games have both SSSE and FPE and that they coincide. The operator formalism is instrumental in this proof. We also address the classes of Zero-Sum Games and Acyclic Games. This analysis is developed in Section 3.

We then illustrate different situations with specific examples. In a first case, an FPE and an SSSE exist and coincide, although the game does not have MFS (the assumption of our main existence result). In a second case, depending on the parameters: either no SSSE exist, or no FPE exist, or both an SSSE and an FPE exist but they do not coincide. In a third case, an FPE exists but Value Iteration does not necessarily converge to it. These examples are summarized in Section 4 and described in more detail in Appendix C.

Finally, we take advantage of the convergence properties we have shown, to propose a solution methodology in a dynamic security game, representing the problem of security patrols in a network. This is reported in Section 5.

1.4 Notation and definitions

We introduce now formally the elements of the model and the notation. A synthesis of this notation is presented in Appendix A.

Let \mathcal{S} represent the finite set of states of the game. Let \mathcal{A}, \mathcal{B} denote the finite set of actions available to players A and B respectively, and we denote by $\mathcal{A}_s \subset \mathcal{A}$ and $\mathcal{B}_s \subset \mathcal{B}$ the actions available in state $s \in \mathcal{S}$. For a given state $s \in \mathcal{S}$ and actions $a \in \mathcal{A}_s$ and $b \in \mathcal{B}_s$, $Q^{ab}(s'|s)$ represents the transition probabilities of reaching the state $s' \in \mathcal{S}$. We denote with Q the family of these probability distributions. The reward received by each player in state s when selecting actions $a \in \mathcal{A}_s$ and $b \in \mathcal{B}_s$ is referred to as the one-step reward functions and are given by $r_A = r_A^{ab}(s)$ and $r_B = r_B^{ab}(s)$. The constants $\beta_A, \beta_B \in [0, 1)$ are discount factors for Player A and B respectively. In our setting time increases discretely and the time horizon is infinite. Therefore we represent a two-person stochastic discrete game \mathcal{G} by

$$\mathcal{G} = (\mathcal{S}, \mathcal{A}, \mathcal{B}, Q, r_A, r_B, \beta_A, \beta_B) .$$

Strategies. We denote by $\mathbb{P}(\mathcal{A}_s)$ and $\mathbb{P}(\mathcal{B}_s)$ the sets of distribution functions over \mathcal{A}_s and \mathcal{B}_s , respectively. The sets of stationary strategies are defined by:

$$\begin{aligned} W_A &= \{f : \mathcal{S} \rightarrow \mathbb{P}(\mathcal{A}) \mid f(s) \in \mathbb{P}(\mathcal{A}_s)\} \\ W_B &= \{g : \mathcal{S} \rightarrow \mathbb{P}(\mathcal{B}) \mid g(s) \in \mathbb{P}(\mathcal{B}_s)\} . \end{aligned}$$

For $f \in W_A$, $f(s)$ is a probability measure on \mathcal{A}_s . In order to simplify the notation, we represent with $f(s, a) = (f(s))(\{a\}) = f(a|s)$ the probability that Player A chooses action a when in state s . Likewise, for $g \in W_B$, we denote with $g(s, b) = g(b|s)$ the probability that Player B chooses b when in state s . In the case that $g \in W_B$ is a deterministic policy,

we will denote directly with $g(s)$ the element of \mathcal{B}_s that has probability one. The notation will be clear from context. The set W_B is assumed to be equipped with a total order \prec_B , which will be used for determining a unique element in case Player B is indifferent between several policies.

In order to simplify notation, given stationary strategies f and g , we define the reward for player $i (= A, B)$ by:

$$r_i^{fg}(s) = \sum_{a \in \mathcal{A}_s} \sum_{b \in \mathcal{B}_s} f(s, a)g(s, b)r_i^{ab}(s). \quad (1.1)$$

Values. Given a pair $(f, g) \in W_A \times W_B$, the evolution of the states is that of a Markov chain on \mathcal{S} with transition probabilities $Q^{fg}(s'|s) = \sum_{a \in \mathcal{A}_s} \sum_{b \in \mathcal{B}_s} f(s, a)g(s, b)Q^{ab}(s'|s)$. Denote with $\{S_n\}_n$ the (random) sequence of states of this Markov chain and \mathbb{E}_s^{fg} the expectation corresponding to the distribution of this sequence, conditioned on the initial state being $S_0 = s$. Then the value of this pair of strategies for Player i , from state s , is:

$$\begin{aligned} V_i^{fg}(s) &= \mathbb{E}_s^{fg} \left[\sum_{k=0}^{\infty} \beta_i^k r_i^{f(S_k), g(S_k)}(S_k) \right] \\ &= \mathbb{E}_s^{fg} \left[\sum_{k=0}^{\infty} \beta_i^k \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} f(S_k, a)g(S_k, b)r_i^{ab}(S_k) \right]. \end{aligned} \quad (1.2)$$

Reaction sets. We proceed with the definition of the player's reaction sets. These definitions rely heavily on the fact that when the leader selects a stationary strategy, the follower faces a finite-state, finite-action, discounted MDP. It is then well-known that there exists optimal stationary and deterministic policies which maximize simultaneously the follower's values starting from any state. Moreover, the set of optimal policies is the cartesian product of the set of optimal decisions in each state. This fact results from *e.g.* Corollary 6.2.8, p. 153 in [13].

Accordingly, let:

$$R_B(f) := \{g \in W_B \mid V_B^{fg}(s) \geq V_B^{fh}(s), \forall s \in \mathcal{S}, h \in W_B\} \cap \prod_{s \in \mathcal{S}} \{0, 1\}^{|\mathcal{B}_s|} \quad (1.3)$$

$$SR_B(f) := \{g \in R_B(f) \mid V_A^{fg}(s) \geq V_A^{fh}(s), \forall s \in \mathcal{S}, h \in R_B(f)\} \quad (1.4)$$

$$\gamma_B(f) := \max_{\prec_B} SR_B(f)$$

$$R_A(s) := \{f \in W_A \mid V_A^{f\gamma_B(f)}(s) \geq V_A^{h\gamma_B(h)}(s), \forall h \in W_A\}. \quad (1.5)$$

Given that Player A selects strategy f , $R_B(f)$ represents the set of deterministic best-response strategies of Player B. As argued above, this set is nonempty. The set $SR_B(f)$ is that of *strong* best-responses, which break ties in favor of Player A. It is possible to break ties simultaneously in all states s , because optimal policies of the MDP form a cartesian product. We denote by $\gamma_B(f)$ the deterministic policy that is the actual best response of Player B to Player A's f . Finally, $R_A(s)$ is the set of Player A's best strategies when starting from state s .

Equilibria. With these notations, we can now define Strong Stackelberg Equilibria of the dynamic game, called here Stationary SSE, as the SSE for the static game where players use stationary strategies in $W_A \times W_B$. It corresponds to the definitions in [11, 16].

Definition 1 (SSSE). A strategy pair $(f, g) \in W_A \times W_B$ is a Stationary Strong Stackelberg Equilibrium if

- i/ $g = \gamma_B(f)$;
- ii/ $f \in R_A(s)$ for all $s \in \mathcal{S}$.

In an SSSE, the strategy f maximizes *simultaneously* the leader’s reward in every state. In contrast with MDP where this is always possible, there is no guarantee that this will happen in a Stackelberg stochastic game. Indeed, in Section 4.3 we provide an example where $\cap_s R_A(s)$ is empty, and consequently there is no SSSE.

To the best of our knowledge, the literature does not provide general statements about the existence of an SSSE. We address in Section 3 this issue in special cases.

2 Operators, Fixed Points and Algorithms

In this section, we develop the formalism of operators, as commonly found in texts on MDPs [13], and also for games in [2, 7, 18]. We focus on fixed points of these operators, as a means to discuss existence of equilibria, and also as a computational procedure. Accordingly, we study the monotonicity and contractivity of these operators. This allows us to prove the convergence of Value Iteration, Policy Iteration and Mathematical Programming-based algorithms, in certain situations.

2.1 Definition of operators

We start with the definition of one-step (or “return function” [7]) operators. The set of value functions, i.e. mappings from \mathcal{S} to \mathbb{R} , will be denoted with $\mathcal{F}(\mathcal{S})$. Given $(f, g) \in W_A \times W_B$ we define $T_i^{fg} : \mathcal{F}(\mathcal{S}) \rightarrow \mathcal{F}(\mathcal{S})$, such that

$$\left(T_i^{fg}v\right)(s) = \sum_{a \in \mathcal{A}_s} f(s, a) \sum_{b \in \mathcal{B}_s} g(s, b) \left[r_i^{ab}(s) + \beta_i \sum_{z \in \mathcal{S}} Q^{ab}(z|s)v(z) \right].$$

It is important to note that the value $(T_i^{fg}v)(s)$ depends only on $f(s)$ and $g(s)$, and not on the rest of the strategies f and g . In the following, with a slight abuse of notation, we will use this quantity for values of f and g specified only at state s .

The set of pairs of value functions is $\mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$. A typical element of it will be denoted as $v = (v_A, v_B)$. Using T_A^{fg} and T_B^{fg} we define the operator T^{fg} on $\mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$ as:

$$(T^{fg}v)_i = T_i^{fg}v_i$$

for $i = A, B$.

It will be recalled in Lemma 1 that T_i^{fg} is a contraction for $i = A, B$. It follows that T^{fg} is contractive as well. As a consequence of Banach's theorem, it admits a unique fixed point on the complete space $\mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$ with the supremum norm, that turns out to be $V^{fg} = (V_A^{fg}, V_B^{fg})$, these functions being defined in (1.2).

Extended reaction sets. We now extend the definitions of reaction sets to involve value functions. They correspond to a dynamic game with only one step and a ‘‘scrap value’’ $v = (v_A, v_B)$. In contrast to the sets introduced in Section 2.1 for SSSE, the sets we discuss here are relative to *local* strategies depending on each state, rather than *global* strategies in W_A and W_B .

For $s \in \mathcal{S}$, $f \in \mathbb{P}(\mathcal{A}_s)$, $v \in \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$, and $v_B \in \mathcal{F}(\mathcal{S})$, let:

$$R_B(s, f, v_B) := \{g \in \mathcal{B}_s \mid (T_B^{fg} v_B)(s) \geq (T_B^{fh} v_B)(s), \forall h \in \mathcal{B}_s\} \quad (2.1)$$

$$SR_B(s, f, v) := \{g \in R_B(s, f, v_B) \mid (T_A^{fg} v_A)(s) \geq (T_A^{fh} v_A)(s), \forall h \in R_B(s, f, v_B)\} \quad (2.2)$$

$$\gamma_B(s, f, v) := \max_{\prec_B} SR_B(s, f, v) \quad (2.3)$$

$$R_A(s, v) := \{f(s) \in \mathbb{P}(\mathcal{A}_s) \mid (T_A^{f\gamma_B(s,f,v)} v_A)(s) \geq (T_A^{h\gamma_B(s,h,v)} v_A)(s), \forall h \in \mathbb{P}(\mathcal{A}_s)\}. \quad (2.4)$$

The definition of Player B's response in (2.3) is such that one unique, non-ambiguous policy is defined as a solution. Any $f(s) \in R_A(s, v)$ is considered as a solution of the problem.

The dynamic programming operator. The one-step Strong Stackelberg problem naturally leads to a mapping in the space of value functions, which is formalized as follows.

Definition 2 (Dynamic programming operator T). Let $T: \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S}) \rightarrow \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$ be defined as:

$$(Tv)_i(s) = \left(T_i^{R_A(s,v), \gamma_B(s, R_A(s,v), v)} v_i \right) (s) \quad (2.5)$$

for $i = A, B$.

Observe that the definition depends on the ordering \prec_B . By changing the ordering, many operators can be defined for the same problem.

Fixed points. We are now in position to define the fixed-point equilibria.

Definition 3 (Fixed Point Equilibrium, FPE). A strategy pair $(f, g) \in W_A \times W_B$ is an FPE if the function $v^* = V^{fg}$, the unique fixed point of T^{fg} , is such that $Tv^* = v^*$. Equivalently, if

$$i/ \quad g(s) = \gamma_B(s, f, v^*) \text{ for all } s \in \mathcal{S};$$

ii/ $f(s) \in R_A(s, v^*)$ for all $s \in \mathcal{S}$.

One of the purposes of this paper is to propose results concerning FPEs and SSSE: discuss whether they respectively exist, and when they do, whether they coincide or not.

2.2 Properties of operators

The following property is well-known (e.g. [13]) for finite-state, finite-action discounted Markov Reward Processes:

Lemma 1. For $i = A, B$, the operator T_i^{fg} is linear, monotone, contractive and V_i^{fg} defined in (1.2) is its unique fixed point. This fixed point has the expression

$$V_i^{fg} = (I - \beta_i Q^{fg})^{-1} r_i^{fg}, \quad (2.6)$$

where r_i^{fg} is defined in (1.1), where the probability transition matrix Q^{fg} is defined similarly in Section 1.4, and I is the identity matrix of appropriate dimension.

We now introduce a particular class of games, and the particular properties of operators for these games.

Definition 4 (Myopic Follower Strategy, MFS). A stochastic game \mathcal{G} is said to be with *Myopic Follower Strategy* if $R_B(s, f, v_B) = R_B(s, f)$, for all $s \in \mathcal{S}$, $f \in W_A$ and $v_B \in \mathcal{F}(\mathcal{S})$.

Games with *Myopic Follower Strategy* are games where the response of the follower is independent of the expected future values: for her, only immediate rewards are relevant. As it is stated in Lemma 5, this setting happens whenever $\beta_B = 0$ (as it is used in [5]) or the leader controls the transitions of the games, see for e.g. single-controller games in [8].

When a game is with MFS, the reaction of the follower depends only on the leader's value v_A :

$$\forall f \in W_A, \forall v \in \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S}), \forall s \in \mathcal{S}, \quad \gamma_B(s, f, v) = \bar{\gamma}_B(s, f, v_A). \quad (2.7)$$

Then the following Lemma is relevant.

Lemma 2. Assume (2.7) holds. Then there exists an operator \bar{T}_A from $\mathcal{F}(\mathcal{S})$ to $\mathcal{F}(\mathcal{S})$ such that for all $v \in \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S})$,

$$(Tv)_A = \bar{T}_A v_A. \quad (2.8)$$

Proof. According to Definition (2.4) and due to (2.7), we have $R_A(s, v) = \bar{R}_A(s, v_A)$ for all $v \in \mathcal{F}(\mathcal{S})$ and $s \in \mathcal{S}$. Then, from (2.5),

$$\begin{aligned} (Tv)_A(s) &= (T_A^{R_A(s,v), \gamma_B(s, R_A(s,v), v)} v_A)(s) \\ &= (T_A^{\bar{R}_A(s, v_A), \bar{\gamma}_B(s, \bar{R}_A(s, v_A), v_A)} v_A)(s) \\ &=: (\bar{T}_A v_A)(s). \end{aligned}$$

□

An alternate construction of operator \bar{T}_A is as follows. It is possible to define, for each $f \in W_A$, the operator \bar{T}_A^f from $\mathcal{F}(\mathcal{S})$ to $\mathcal{F}(\mathcal{S})$ as:

$$(\bar{T}_A^f v_A)(s) = (T_A^{f, \bar{\gamma}_B(s, f, v_A)} v_A)(s) . \quad (2.9)$$

Another consequence of MFS is this property which follows from the definition of $\gamma_B(\cdot)$ in (2.3) and $SR_B(\cdot)$ in (2.2):

$$(\bar{T}_A^f v_A)(s) = \max_{g \in R_B(s, f)} (T_A^{fg} v_A)(s) . \quad (2.10)$$

In this equation, the maximization set on the right-hand side does not depend on v at all. Finally, define the operator \bar{T}_A from $\mathcal{F}(\mathcal{S})$ to $\mathcal{F}(\mathcal{S})$ as, for all $s \in \mathcal{S}$:

$$(\bar{T}_A v_A)(s) = \max_{f \in W_A} (\bar{T}_A^f v_A)(s) . \quad (2.11)$$

Observe that the maximum is indeed attained, because the right-hand side is a linear combination of the finite set of values $f(s, a)$, $a \in \mathcal{A}_s$.

We can now state the principal tool of this paper for ascertaining the existence of FPE.

Theorem 1. Let \mathcal{G} be a stochastic game with MFS, then it is true that:

- a) For any stationary strategy $f \in W_A$, the operator $\bar{T}_A^f : \mathcal{F}(\mathcal{S}) \rightarrow \mathcal{F}(\mathcal{S})$, defined in (2.9) is a contraction on $(\mathcal{F}(\mathcal{S}), \|\cdot\|_\infty)$ of modulus β_A .
- b) The operator \bar{T}_A defined in (2.11) is a contraction on $(\mathcal{F}(\mathcal{S}), \|\cdot\|_\infty)$, of modulus β_A .
- c) For any stationary strategy $f \in W_A$, operator \bar{T}_A^f is monotone. Operator \bar{T}_A is monotone as well.

Proof. The central argument of the proof is the following fact. Let h_1 and h_2 be two real functions defined on some set B , where they attain their maximum. Then for all $b_1 \in \arg \max_B \{h_1(b)\}$ and all $b_2 \in \arg \max_B \{h_2(b)\}$,

$$h_1(b_2) - h_2(b_2) \leq \max_B \{h_1(b)\} - \max_B \{h_2(b)\} \leq h_1(b_1) - h_2(b_1) . \quad (2.12)$$

To show a), take $v_A, u_A \in \mathcal{F}(\mathcal{S})$, a stationary strategy $f \in W_A$, and $s \in \mathcal{S}$. Then, using (2.10),

$$\begin{aligned} (\bar{T}_A^f v_A)(s) - (\bar{T}_A^f u_A)(s) &= \max_{b \in R_B(s, f)} \sum_{a \in \mathcal{A}_s} f(s, a) \left[r_A^{ab}(s) + \beta_A \sum_{z \in \mathcal{S}} Q^{ab}(z|s) v_A(z) \right] \\ &\quad - \max_{b \in R_B(s, f)} \sum_{a \in \mathcal{A}_s} f(s, a) \left[r_A^{ab}(s) + \beta_A \sum_{z \in \mathcal{S}} Q^{ab}(z|s) u_A(z) \right] . \end{aligned} \quad (2.13)$$

Then from (2.12), there exists $b \in R_B(s, f)$ such that:

$$\begin{aligned}
(\bar{T}_A^f v_A)(s) - (\bar{T}_A^f u_A)(s) &\leq \sum_{a \in \mathcal{A}_s} f(s, a) \left[r_A^{ab}(s) - r_A^{ab}(s) \right. \\
&\quad \left. + \beta_A \sum_{z \in \mathcal{S}} (Q^{ab}(z|s)v_A(z) - Q^{ab}(z|s)u_A(z)) \right] \\
&= \sum_{a \in \mathcal{A}_s} f(s, a) \beta_A \sum_{z \in \mathcal{S}} Q^{ab}(z|s)(v_A(z) - u_A(z)) \quad (2.14) \\
&\leq \beta_A \|v_A - u_A\|_\infty .
\end{aligned}$$

By reversing the roles of v_A and u_A , then taking the maximum over $s \in \mathcal{S}$ we have that:

$$\|\bar{T}_A^f v_A - \bar{T}_A^f u_A\|_\infty = \max_{s \in \mathcal{S}} |(\bar{T}_A^f v_A)(s) - (\bar{T}_A^f u_A)(s)| \leq \beta_A \|v_A - u_A\|_\infty ,$$

concluding that \bar{T}_A^f is a contracting map of modulus β_A .

In order to show b), take $v_A, u_A \in \mathcal{F}(\mathcal{S})$, a state $s \in \mathcal{S}$ and f^* any optimal policy realizing $\max_{f \in W_A} (\bar{T}_A^f v_A)(s)$. Then,

$$\begin{aligned}
(\bar{T}_A v_A)(s) - (\bar{T}_A u_A)(s) &= \max_{f \in \mathbb{P}(A_s)} (\bar{T}_A^f v_A)(s) - \max_{f \in \mathbb{P}(A_s)} (\bar{T}_A^f u_A)(s) \\
&= (\bar{T}_A^{f^*} v_A)(s) - \max_{f \in \mathbb{P}(A_s)} (\bar{T}_A^f u_A)(s) \\
&\leq (\bar{T}_A^{f^*} v_A)(s) - (\bar{T}_A^{f^*} u_A)(s) \\
&\leq \beta_A \|v_A - u_A\|_\infty .
\end{aligned}$$

Then, by reversing the roles of v_A, u_A and taking the maximum the result follows.

Consider now statement c). If $v_A \geq u_A$, then from (2.13) with (2.12) and (2.14), there exists some $b \in R_B(s, f)$ such that:

$$\begin{aligned}
(\bar{T}_A^f v_A)(s) - (\bar{T}_A^f u_A)(s) &\geq \sum_{a \in \mathcal{A}_s} f(s, a) \beta_A \sum_{z \in \mathcal{S}} Q^{ab}(z|s)(v_A(z) - u_A(z)) \\
&\geq 0 .
\end{aligned}$$

Then c) also holds. The extension to \bar{T}_A is classical. \square

2.3 Value Iteration Algorithms

Value Iteration generally consists in applying a dynamic programming operator to some initial value function, until a convergence criterion is met. Specifically, given some $\varepsilon > 0$, Value Iteration applies some operator repeatedly until the distance between two functions v_A^n and v_A^{n+1} is less than ε .

In view of the preceding discussion, two variants of the algorithm will be used: one for the general situation (Algorithm 1) and one for the specific situation of MFS, i.e. when (2.7) holds (Algorithm 2).

Algorithm 1 Value function iteration for infinite horizon; general case

Require: $\varepsilon > 0$

- 1: Initialize with $n = 0$, $v_A^0(s) = v_B^0(s) = 0$ for every $s \in \mathcal{S}$
 - 2: **repeat**
 - 3: $n := n + 1$
 - 4: Compute v^n as $v^n(s) := (Tv^{n-1})(s)$, $\forall s \in \mathcal{S}$
with T according to Definition (2.5)
 - 5: **until** $\|v^n - v^{n-1}\|_\infty \leq \varepsilon$
 - 6: Pick (f^*, g^*) such that $v^n(s) = (T^{f^*g^*}v^{n-1})(s)$ for all $s \in \mathcal{S}$
 - 7: **return** Approximate Stationary Strong Stackelberg policies (f^*, g^*)
-

Algorithm 2 Value function iteration for infinite horizon; simplified case

Require: $\varepsilon > 0$

- 1: Initialize with $n = 0$, $v_A^0(s) = 0$ for every $s \in \mathcal{S}$
 - 2: **repeat**
 - 3: $n := n + 1$
 - 4: Compute v_A^n as $v_A^n(s) := (\bar{T}_A v_A^{n-1})(s)$, $\forall s \in \mathcal{S}$
with \bar{T} as in (2.11)
 - 5: **until** $\|v_A^n - v_A^{n-1}\|_\infty \leq \varepsilon$
 - 6: Pick f^* such that $v_A^n(s) = (\bar{T}_A^{f^*} v_A^{n-1})(s)$ and g^* such that $g^*(s) = \gamma_B(s, f^*, v_A^{n-1})$, for all $s \in \mathcal{S}$
 - 7: **return** Approximate Stationary Strong Stackelberg policies (f^*, g^*)
-

There is no guarantee in general that Algorithm 1 will converge, and we present in Section 4.4 an example where it does not. However, thanks to Theorem 2, we can state that Algorithm 2 does converge.

Theorem 2. Let \mathcal{G} be a stochastic game with MFS. Then the sequence of value functions v_A^n in Algorithm 2 converges to v_A^* , which is the fixed point of \bar{T}_A . Moreover the following bounds hold:

$$\|v_A^* - v_A^n\|_\infty \leq \frac{2\beta_A^n \|r_A\|_\infty}{1 - \beta_A} \quad \text{for any } n \in \mathbb{N},$$

$$\|v_A^* - V_A^{f^*g^*}\|_\infty \leq \frac{2\beta_A \varepsilon}{1 - \beta_A}.$$

Proof. Let the pair of policies (f^*, g^*) be the ones returned by Algorithm 2 and $V_A^{f^*g^*}$ be the fixed point of $\bar{T}_A^{f^*}$. By Theorem 1 b) and Banach's Theorem, we know that \bar{T}_A has a unique fixed point, v_A^* . Then,

$$\|V_A^{f^*g^*} - v_A^*\|_\infty \leq \|V_A^{f^*g^*} - v_A^n\|_\infty + \|v_A^n - v_A^*\|_\infty. \quad (2.15)$$

In (2.15), the first term on the right-hand side is bounded as follows:

$$\begin{aligned}
\|V_A^{f^*g^*} - v_A^n\|_\infty &= \|\bar{T}_A^{f^*} V_A^{f^*g^*} - v_A^n\|_\infty \\
&\leq \|\bar{T}_A^{f^*} V_A^{f^*g^*} - \bar{T}_A v_A^n\|_\infty + \|\bar{T}_A v_A^n - v_A^n\|_\infty \\
&= \|\bar{T}_A^{f^*} V_A^{f^*g^*} - \bar{T}_A^{f^*} v_A^n\|_\infty + \|\bar{T}_A v_A^n - \bar{T}_A v_A^{n-1}\|_\infty \\
&\leq \beta_A \|V_A^{f^*g^*} - v_A^n\|_\infty + \beta_A \|v_A^n - v_A^{n-1}\|_\infty,
\end{aligned}$$

where the first equality is by definition, the inequality right after is the triangular inequality. The third line is because of the definition of $\bar{T}_A^{f^*}$ and the inequality is because $\bar{T}_A^{f^*}$ and \bar{T}_A are contracting maps of modulus β_A . This last inequality implies that

$$\|V_A^{f^*g^*} - v_A^n\|_\infty \leq \frac{\beta_A}{1 - \beta_A} \|v_A^n - v_A^{n-1}\|_\infty.$$

For the second term on the right-hand side of (2.15), we have:

$$\begin{aligned}
\|v_A^n - v_A^*\|_\infty &= \lim_{t \rightarrow \infty} \|v_A^n - v_A^t\|_\infty \\
&\leq \lim_{t \rightarrow \infty} \sum_{k=0}^{t-n-1} \|v_A^{n+k} - v_A^{n+k+1}\|_\infty \\
&\leq \lim_{t \rightarrow \infty} \sum_{k=0}^{t-n-1} \beta_A^k \|v_A^{n-1} - v_A^n\|_\infty \\
&= \frac{\beta_A}{1 - \beta_A} \|v_A^{n-1} - v_A^n\|_\infty.
\end{aligned}$$

Then, we have that the policies returned by the algorithm satisfy:

$$\|V_A^{f^*g^*} - v_A^*\|_\infty \leq 2 \frac{\beta_A}{1 - \beta_A} \|v_A^{n-1} - v_A^n\|_\infty = \frac{2\beta_A}{1 - \beta_A} \varepsilon.$$

Furthermore, given that

$$\|v_A^{n-1} - v_A^n\|_\infty \leq \beta_A^{n-1} \|v_A^0 - v_A^1\|_\infty = \beta_A^{n-1} \|v_A^1\|_\infty \leq \beta_A^{n-1} \|r_A\|_\infty$$

the result follows. \square

2.4 Policy Iteration

The Policy Iteration (PI) algorithm directly iterates in the policy space. This algorithm starts with an arbitrary policy f and then finds the optimal infinite discounted horizon values, taking into account the optimal response $g(f)$. These values are then used to compute new policies. These two steps of the algorithm can be defined as the *Evaluation Phase* and the *Improvement Phase*.

As in the previous section, two variants of the algorithm will be used: one for the general situation (Algorithm 3) and one for the specific situation of a MFS, i.e. when (2.7) holds (Algorithm 4).

Algorithm 3 Policy Iteration (PI); general case

- 1: Require $\varepsilon > 0$
 - 2: Initialize with $n = 0$
 - 3: Choose an arbitrary pair of strategies $(f_0, g_0) \in W_A \times W_B$ with $g_0(s) = \gamma_B(s, f_0, \mathbf{0})$ for all $s \in \mathcal{S}$
 - 4: Compute $u^0 = (u_A^0, u_B^0)$ fixed point of $T^{f_0 g_0}$
 - 5: **repeat**
 - 6: $n := n + 1$
 - 7: **Improvement Phase:** Find a pair of strategies (f_n, g_n) such that $T^{f_n g_n} u^{n-1} = T u^{n-1}$ with $g_n(s) = \gamma_B(s, f_n, u^{n-1})$ for all $s \in \mathcal{S}$
 - 8: **Evaluation Phase:** Find $u^n = (u_A^n, u_B^n)$, fixed point of the operator $T^{f_n g_n}$
 - 9: **until** $\|u^n - u^{n-1}\|_\infty \leq \varepsilon$
 - 10: $f^* := f_n; g^*(s) := \gamma_B(s, f_n, u^n)$ for all $s \in \mathcal{S}$
 - 11: **return** Approximate Stationary Strong Stackelberg policies (f^*, g^*)
-

Algorithm 4 Policy Iteration (PI); simplified case

- 1: Require $\varepsilon > 0$
 - 2: Initialize with $n = 0$
 - 3: Choose an arbitrary pair of strategies $(f_0, g_0) \in W_A \times W_B$ with $g_0(s) = \bar{\gamma}_B(s, f_0, \mathbf{0})$ for all $s \in \mathcal{S}$
 - 4: Compute u_A^0 fixed point of $\bar{T}_A^{f_0}$
 - 5: **repeat**
 - 6: $n := n + 1$
 - 7: **Improvement Phase:** Find a distribution f_n such that $\bar{T}_A^{f_n} u_A^{n-1} = \bar{T}_A u_A^{n-1}$
 - 8: **Evaluation Phase:** Find u_A^n fixed point of the operator $\bar{T}_A^{f_n}$
 - 9: **until** $\|u_A^n - u_A^{n-1}\|_\infty \leq \varepsilon$
 - 10: $f^* := f_n; g^*(s) := \bar{\gamma}_B(s, f_n, u_A^n)$ for all $s \in \mathcal{S}$
 - 11: **return** Approximate Stationary Strong Stackelberg policies (f^*, g^*)
-

The Evaluation Phase in Algorithm 3 (respectively Algorithm 4) requires to solve two (resp. one) linear systems of size $|\mathcal{S}| \times |\mathcal{S}|$. On the other hand, the Improvement Phase can be implemented by solving a static Strong Stackelberg equilibrium for each state $s \in \mathcal{S}$. Now we prove that Algorithm 4 converges to the SSSE. In other words, the PI algorithm converges to the SSSE for stochastic games with MFS.

Lemma 3. If a function $v_A \in \mathcal{F}(\mathcal{S})$ satisfies $v_A \leq \bar{T}_A^f v_A$, for some $f \in \mathbb{P}(\mathcal{A})$ then $v_A \leq v_A^f$, where v_A^f is the unique fixed point of \bar{T}_A^f in $\mathcal{F}(\mathcal{S})$. The same holds for \bar{T}_A .

Proof. By hypothesis we have that

$$v_A \leq \bar{T}_A^f v_A ,$$

that implies by Theorem 1 c),

$$\bar{T}_A^f v_A \leq (\bar{T}_A^f)^2 v_A ,$$

and then

$$v_A \leq (\bar{T}_A^f)^2 v_A .$$

In the same way, for each n we have

$$v_A \leq (\bar{T}_A^f)^n v_A ,$$

and by Theorem 1 a), when $n \rightarrow \infty$,

$$(\bar{T}_A^f)^n v_A \longrightarrow v_A^f .$$

The result follows for \bar{T}_A^f . The proof for \bar{T}_A is similar. \square

Theorem 3. Suppose that Condition (2.7) holds. The sequence of functions u_A^n in Algorithm 3 verifies $u_A^n \uparrow v_A^*$. Further, if for any $n \in \mathbb{N}$, $u_A^n = u_A^{n+1}$, then it is true that $u_A^n = v_A^*$.

Proof. For each $s \in \mathcal{S}$, we have that

$$u_A^0(s) = \bar{T}_A^{f_0}(u_A^0)(s) \leq \bar{T}_A(u_A^0)(s) = \bar{T}_A^{f_1}(u_A^0)(s) .$$

Then the value function u_A^0 satisfies

$$u_A^0 \leq \bar{T}_A^{f_1} u_A^0 ,$$

and by Lemma 3

$$u_A^0 \leq v_A^{f_1} = u_A^1 .$$

Iterating over n , we have that

$$u_A^n \leq \bar{T}_A(u_A^n) \leq u_A^{n+1} . \quad (2.16)$$

Now the sequence $\{u_A^n\}_{n \in \mathbb{N}}$ being non-decreasing and bounded by $\|r_A\|_\infty / (1 - \beta_A)$, there exists a value function u_A such that for any $s \in \mathcal{S}$

$$u_A(s) = \lim_{n \rightarrow \infty} u_A^n(s) .$$

Taking $n \rightarrow \infty$ in (2.16), $u_A \leq \bar{T}_A(u_A) \leq u_A$ and therefore $u_A = \bar{T}_A(u_A)$, and by uniqueness of the fixed point

$$u_A = v_A^* ,$$

and we have the first claim of the theorem: $u_A^n \uparrow v_A^*$. Also, if it is verified for some n that $u_A^n = u_A^{n+1}$, then, using (2.16),

$$u_A^{n+1} = u_A^n \leq \bar{T}_A u_A^n \leq u_A^{n+1},$$

which implies

$$u_A^n = \bar{T}_A u_A^n = v_A^*,$$

where the second equality is again given by the uniqueness of the fixed point. The second claim follows. \square

The results exposed in this section strongly rely on the fact that $\gamma_B(s, f, v)$ is independent on v_B . In Section 3 we show that MFS is a sufficient condition for the existence of an FPE but all the results here may fail in the general case.

2.5 Mathematical Programming Formulations

In this section we develop the discussion of Mathematical Programming (MP) formulations, as the one proposed in [16]. To start the discussion we notice that for each $f \in W_A$ the follower solves an MDP with transition and rewards given by the expectation induced by f . Then, as argued in Section 1.4, there exists (at least) one optimal policy in the set of deterministic stationary policies. This policy g can be retrieved by finding deterministic policies that induce a fixed point of the operator $T_B^{f\gamma(s, f, v)}$. This condition is modeled as the following set of non-linear constraints:

$$0 \leq v_B(s) - (T_B^{fg} v_B)(s) \leq M_B(1 - g_{sb}) \quad s \in \mathcal{S}, b \in \mathcal{B}_s \quad (2.17)$$

$$\sum_{b \in \mathcal{B}_s} g_{sb} = 1 \quad s \in \mathcal{S} \quad (2.18)$$

$$g_{sb} \in \{0, 1\} \quad s \in \mathcal{S}, b \in \mathcal{B}_s. \quad (2.19)$$

The variable of the program g_{sb} is meant to represent the probability $g(s, b)$.

For each $f \in W_A$, the deterministic best response set of the follower is determined by constraints (2.17)–(2.19). In (2.17) the constant M_B is chosen so that when $g_{sb} = 0$, the upper bound is not constraining. Since $\|v_B\|_\infty \leq \|r_B\|_\infty / (1 - \beta_B)$, the value $M_B = 2\|r_B\|_\infty / (1 - \beta_B)$ is adequate. Now the leader's problem can be reduced to determine which f maximizes in each state its total expected reward. Vorobeychik and Singh in [16] propose the following formulation:

$$(MP) \quad \max \sum_{s \in \mathcal{S}} \alpha_s v_A(s) \quad (2.20)$$

s.t. Constraints (2.17), (2.18), (2.19)

$$\sum_{a \in \mathcal{A}_s} f_{sa} = 1 \quad s \in \mathcal{S}$$

$$v_A(s) - (T_A^{fg} v_A)(s) \leq M_A(1 - g_{sb}) \quad s \in \mathcal{S}, b \in \mathcal{B}_s \quad (2.21)$$

$$f_{sa} \geq 0, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}_s.$$

where $\alpha \in \mathbb{R}_+^{|\mathcal{S}|}$ is a non-negative vector of coefficients. Problem (MP) above is a non-linear optimization problem with integer variables. Such problems are challenging to solve in general. In particular, constraints (2.17) and (2.21) are non-convex quadratic constraints that involve integer variables.

This optimization problem is built based on an analogy to MDPs. In particular, it uses the reduction to a single objective using a vector of weights in (2.20). For MDPs, this choice is arbitrary. As it turns out in the experiments presented in Section 4, the result of (MP) sometimes depends on the vector α_s , and sometimes it is not an SSSE. We observe that when such anomalies occur, the operator T_A is not monotone. On the other hand, in cases where T_A is known to be monotone, no such problems seem to occur. We therefore conjecture that for the correctness of (MP), it is necessary to have monotonicity of the operator.

3 Existence Results for SSSE and FPE.

In this section we present existence results for Stationary Strong Stackelberg Equilibria and Fixed-Points Equilibria of the dynamic programming operator. The general idea underlying these results is that, under certain assumptions, it can be proved that some operator, typically T or \bar{T}_A defined in Section 2, is contractive. Using Banach's theorem, it then has a fixed point with which the solution is constructed. Then, still under some assumptions, this FPE solution is shown to be an SSSE.

3.1 Single-state results

In the case where there is only one state, the game is equivalent to a static game, so that SSE and SSSE coincide. Indeed, if $\mathcal{S} = \{s_0\}$, it is clear that

$$V_i^{fg}(s_0) = \frac{r_i^{fg}(s_0)}{1 - \beta_i}. \quad (3.1)$$

for all $(f, g) \in W_A \times W_B$, so that optimization of V_i^{fg} and of r_i^{fg} are equivalent.

The existence of an SSSE for single-state games is well accepted in the literature with however no clear reference. In this section, we state and prove this result and connect it to the FPE.

We start with a general and useful result that applies to any game.

Lemma 4. Let \mathcal{G} be a stochastic game. For all $s \in \mathcal{S}$, the set $R_A(s)$ is nonempty.

Proof. We use the scheme of proof of Proposition 3.1 in [15]. Fix a state $s \in \mathcal{S}$. Define $D = \{(f, g) \in W_A \times W_B | g \in R_B(f)\}$. The value functions V_i^{fg} can be expressed as $V_i^{fg} = (I - \beta_i Q^{fg})^{-1} r_i^{fg}$. Due to the finiteness of \mathcal{S} , this is a rational function of f and g . It does not have singularities inside $W_A \times W_B$ and is therefore continuous. In particular, the mappings $(f, g) \mapsto V_i^{fg}(s)$ are continuous. Since W_A and W_B are compact,

the maximum theorem applies: the maximum of this function over D exists. Therefore the set $R_A(s)$ is nonempty. \square

Theorem 4. If the game \mathcal{G} has only one state, then it has an SSSE which is also an FPE.

Proof. Let \mathcal{S} have a single state: $\mathcal{S} = \{s_0\}$. The existence of SSSE is a particular case of Lemma 4. The existence of an FPE follows from the observation that the game is MFS: Theorem 1 applies to it. Being a contraction, the operator \bar{T}_A has a unique fixed point from which an FPE is constructed. There remains to show that this FPE coincides with the SSSE.

To that end, we first show that $R_B(f) = R_B(s_0, f)$ for all $f \in W_A$ (since the game is with MFS, this latter set does not depend on $v_B \in \mathcal{F}(\mathcal{S})$). If $g \in R_B(s_0, f)$, then for all $h \in W_B$,

$$V_B^{fg}(s_0) = \frac{r_B^{fg}(s_0)}{1 - \beta_B} \geq \frac{r_B^{fh}(s_0)}{1 - \beta_B} = V_B^{fh}(s_0)$$

which means that $g \in R_B(f)$. By the same token, if $g \in R_B(f)$ then $g \in R_B(s_0, f)$. So both reaction sets coincide. Since V_A^{fg} and r_A^{fg} are also proportional, breaking ties in favor of the leader is the same problem for SSSE and FPE: the sets $SR_B(f)$ and $SR_B(s_0, f, v_B)$ also coincide, and $\gamma_B(f) = \gamma_B(s_0, f, v)$ for all $f \in W_A$ and $v \in \mathcal{F}(\mathcal{S})$. It follows from (1.5) and (2.4) that $R_A(s_0) = R_A(s_0, v)$ for all v , which means that SSSE and FPE coincide. \square

An alternative algorithmic proof of the existence of the SSSE in Theorem 4 is provided in Appendix B. This is based in the Multiple LPs algorithm which also give as a polynomial algorithm to solve an SSSE in the static case.

3.2 Myopic Follower Strategies

Theorem 5 (FPE for MFS). If the game \mathcal{G} is with MFS then it admits an FPE.

Proof. According to Theorem 1 b), the operator \bar{T}_A is contractive. It therefore admits a fixed point v_A^* . Let $f^* \in W_A$ be defined by, for each $s \in \mathcal{S}$, $f^*(s) = \bar{R}_A(s, v_A^*)$. Let $g^* \in W_B$ be defined for each $s \in \mathcal{S}$, by $g^*(s) = \bar{\gamma}_B(s, f^*, v_A^*) = \gamma_B(s, f^*, v^*)$. We show that (f^*, g^*) is an FPE.

To avoid confusion in the notation, denote with $U = V^{f^*g^*}$, the unique fixed point of $T^{f^*g^*}$. We first check that $v_A^* = U_A$. We have successively: for every $s \in \mathcal{S}$,

$$\begin{aligned} v_A^*(s) &= (\bar{T}_A v_A^*)(s) \\ &= (T_A^{R_A(s, v_A^*), \bar{\gamma}_B(s, R_A(s, v_A^*), v_A^*)} v_A^*)(s) \\ &= (T_A^{f^*(s)g^*(s)} v_A^*)(s) . \end{aligned}$$

The first line is the definition of v_A^* as a fixed point. The second one is the definition of operator \bar{T}_A in (2.8) and that of T in (2.5), combined with the MFS property, see the proof of Lemma 2. The third one is by definition of f^* and g^* . This last line is equivalent to saying that v_A^* is the fixed-point of operator $T_A^{f^*g^*}$, hence $v_A^* = U_A$. As a consequence, $(TU)_A = U_A$.

There remains to be seen that $(TU)_B = U_B$. We have: $(TU)_B = T_B^{f^*g^*} U_B = U_B$ since by definition of U , U_B is the fixed point of $T^{f^*g^*}$. This completes the proof. \square

In the following Lemma 5, we show that there are actually two main classes games which have MFS. We introduce now these classes of games with one important subclass.

Myopic follower: We define a game as a myopic follower game if $\beta_B = 0$. Note that in this case the follower at any step of the game does not take into account the future rewards, but only the instantaneous rewards.

In this case, the one-step operator of the follower is: $(T_B^{fg} v_B)(s) = r_B^{fg}(s)$ (see (1.1)) and it clearly does not depend on v_B . Therefore, the reaction set $R_B(s, f, v_B)$ defined in (1.3) does not depend either on v_B : $R_B(s, f, v_B) = R_B(s, f)$. It follows that the follower's best response has the form (2.7).

Leader-Controller Discounted Games: This case is a particular case of the Single-controller discounted game described in Filar and Vrieze [8], where the controller is the leader. In other words, the transition law has the form $Q^{ab}(z|s) = Q^a(z|s)$.

In that case, the one-step operator of the follower is:

$$\begin{aligned} (T_B^{fg} v_B)(s) &= \sum_{a \in \mathcal{A}_s} f(s, a) \sum_{b \in \mathcal{B}_s} g(s, b) \left[r_B^{ab}(s) + \beta_B \sum_{z \in \mathcal{S}} Q^a(z|s) v_B(z) \right] \\ &= r_B^{fg}(s) + \beta_B \sum_{a \in \mathcal{A}_s} f(s, a) \sum_{z \in \mathcal{S}} Q^a(z|s) v_B(z). \end{aligned}$$

Then, for $g, h \in W_B$, we have: $(T_B^{fg} v_B)(s) - (T_B^{fh} v_B)(s) = r_B^{fg}(s) - r_B^{fh}(s)$ and the difference does not depend on v_B . The reaction set $R_B(s, f, v_B)$ is defined as those g such that: $\forall h \in \mathcal{B}_s$, $(T_B^{fg} v_B)(s) - (T_B^{fh} v_B)(s) \geq 0$. It is therefore independent from v_B for any $s \in \mathcal{S}$, and as before, (2.7) holds.

Multi-stage games: in such games, the state evolves sequentially and deterministically through s_1, s_2, \dots, s_K and stops. This can be seen a particular case of Leader-Controlled Discounted Game, where the evolution is actually not controlled at all. An additional terminal state with trivial reward functions may be needed to model the end of a game with finitely many stages.

The reduction of MFS to these classes is the topic of the following lemma.

Lemma 5. Let \mathcal{G} be a game with MFS. Then one the following statements is true:

i/ $\beta_B = 0$;

ii/ $Q^{ab}(z|s) = Q^a(z|s)$ for all $s, z \in \mathcal{S}$ and all $a \in \mathcal{A}_s, b \in \mathcal{B}_s$.

Proof. We prove by contradiction the following statement:

$$\forall s \in \mathcal{S}, \forall a \in \mathcal{A}_s, \forall b \in \mathcal{B}_s, \forall z \in \mathcal{S}, \quad \beta_B (Q^{ab}(z|s) - Q^{ab'}(z|s)) = 0, \quad (3.2)$$

which itself is equivalent to the statement of the lemma.

For each $a \in \mathcal{A}_s$ and $s \in \mathcal{S}$, consider the policy where the leader plays the pure strategy a , denoted by δ_a . Take $b^* \in R_B(s, \delta_a, v_B)$ for a given v_B (note that $R_B(s, f, v_B) \neq \emptyset$). Then it is true that for all $b \in \mathcal{B}_s$:

$$r_B^{ab^*}(s) - r_B^{ab}(s) + \sum_{z \in \mathcal{S}} \beta_B (Q^{ab^*}(z|s) - Q^{ab}(z|s)) v_B(z) \geq 0.$$

Suppose by contradiction that (3.2) does not hold. Then there exists s, a, b and b' such that $\xi = \beta_B(Q^{ab^*}(z^*|s) - Q^{ab'}(z^*|s)) \neq 0$ for some z^* . Then by taking $v'_B(z^*)$ with the opposite sign of ξ big enough, and $v'_B(z) = 0$, for $z \neq z^*$, the inequality will turn negative. That would mean that b^* does not belong to $R_B(s, f, v'_B)$ with v'_B and then the game is not MFS. This is a contradiction, so such elements s, a, b, b', z do not exist, and (3.2) holds. \square

We now state the principal results of this section: the MFS property implies the existence of both FPE and SSSE, and their coincidence.

Theorem 6. Let \mathcal{G} be a stochastic game with MFS. Then \mathcal{G} has an SSSE, which corresponds to its FPE.

Proof. Let (f^*, g^*) be the FPE of game \mathcal{G} , and $V^* = V^{f^*g^*}$. We know that the FPE exists by Theorem 5. From the proof of this result, we know that $V_A^* = \bar{T}_A V_A^* = \bar{T}_A^{f^*} V_A^*$.

We first prove that $R_B(f^*) = \prod_{s \in \mathcal{S}} R_B(s, f^*)$. According to Lemma 5, since the game has MFS, then either $\beta_B = 0$, or the game is Leader-Controlled Discounted. In both cases, the value of Player B has the form (see (2.6)):

$$V_B^{fg} = (I - \beta_B Q^f)^{-1} r_B^{fg},$$

where Q^f is the leader-controlled transition matrix, relevant only in case $\beta_B \neq 0$. We note that, given that the matrix $(I - \beta_B Q^f)^{-1}$ is positive, $r_B^{fg} \geq r_B^{fh}$, implies $V_B^{fg} \geq V_B^{fh}$. Additionally, if for some s, g, h , $r_B^{fg}(s) > r_B^{fh}(s)$, then $V_B^{fg}(s) > V_B^{fh}(s)$.

Let f be an arbitrary element of W_A . On the one hand, $\prod_s R_B(s, f) \subset R_B(f)$. To see this, pick $g \in \prod_s R_B(s, f)$. Then for all $h \in W_B$, $r_B^{fg} \geq r_B^{fh}$ and therefore $V_B^{fg} \geq V_B^{fh}$: this means $g \in R_B(f)$. The set $R_B(f)$ is therefore nonempty. On the other hand, $R_B(f) \subset \prod_s R_B(s, f)$. To see this, pick $g \in R_B(f)$ (the set is not empty). If it is not in $\prod_s R_B(s, f)$, then there is some s and some $b \in R_B(s, f)$ such that $r_B^{fb}(s) > r_B^{fg}(s)$. Then the policy $h \in W_B$ which coincides with g except at state s where $h(s) = b$, is such that $V_B^{fh}(s) > V_B^{fg}(s)$, a contradiction. Therefore, $\prod_s R_B(s, f) = R_B(f)$, for all $f \in W_A$.

At this point, we have shown that Player B reacts the same way to Player A's strategy f , in the SSSE problem or in the FPE problem with any scrap value function v . Nevertheless, we cannot conclude that the *strong* reaction is the same, since that of the FPE problem *does* depend on the scrap value v .

However, we know that Player B's tie-breaking problem in (1.4) is a Markov Decision Problem. This means that the value of Player A after Player B's strong reaction, say V_A^f , is given by a Bellman equation, namely:

$$\begin{aligned} V_A^f(s) &= \max_{g \in R_B(f)} \{r_A^{fg}(s) + \beta_A(Q^{fg}V_A^f)(s)\} \\ &= \max_{b \in R_B(s,f)} \{r_A^{fb}(s) + \beta_A(Q^{fb}V_A^f)(s)\} \end{aligned} \quad (3.3)$$

for all $s \in \mathcal{S}$. Here, we have used the fact that $R_B(f)$ is a cartesian product, and that MDPs can be solved state by state. We recognize in the right-hand side of (3.3) the operator \bar{T}_A^f defined in (2.10). In other words, V_A^f is the fixed point of \bar{T}_A^f .

Let then define $U_A \in \mathcal{F}(\mathcal{S})$ as:

$$U_A(s) = \max_{f \in W_A} V_A^f(s) = V_A^{f_s}(s).$$

Here, $f_s \in W_A$ realizes the maximum for state s . By construction, $U_A(s) \geq V_A^f(s)$ for any particular $f \in W_A$. We proceed to prove that $U_A = V_A^*$. First, consider the action of \bar{T}_A on U_A : for $s \in \mathcal{S}$,

$$(\bar{T}_A U_A)(s) = \max_{f \in W_A} (\bar{T}_A^f U_A)(s) \geq (\bar{T}_A^{f_s} U_A)(s) \geq (\bar{T}_A^{f_s} V_A^{f_s})(s) = V_A^{f_s}(s) = U_A(s).$$

The first equality is the definition of \bar{T}_A . The first inequality is clear. The second one results from the monotonicity of operator \bar{T}_A^f . The second equality is because $V_A^{f_s}$ is the fixed point of $\bar{T}_A^{f_s}$. Then according to Lemma 3, $\bar{T}_A U_A \geq U_A$ implies $U_A \leq V_A^*$ since V_A^* is the fixed point of \bar{T}_A .

Now, since $V_A^* = V_A^{f^*}$, the fixed point of operator $\bar{T}_A^{f^*}$, then for all $s \in \mathcal{S}$, $U_A(s) = \max_f V_A^f(s) \geq V_A^{f^*}(s) = V_A^*(s)$. In other words, $U_A \geq V_A^*$. We conclude that indeed $U_A = V_A^*$.

As a consequence, we have shown that $f^* \in \bigcap_{s \in \mathcal{S}} R_A(s)$, that is, (f^*, g^*) is an SSSE. \square

3.3 Zero-Sum Games

In zero-sum games, $\beta_A = \beta_B$ and $r_B = -r_A$.

Theorem 7. If the game \mathcal{G} is a Zero-Sum Game, then it admits an FPE.

The existence of an FPE follows from the contractivity of the operator associated, in a similar way as in [7, Section 8] for Nash Equilibria in Stochastic Games. We include here an argument in the line of the proof of Theorem 1.

Proof. Consider a function v in the set $\mathcal{W} = \{v \in \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S}) \mid v_B = -v_A\}$. Since v_B can be substituted with $-v_A$, it turns out that $SR_B(s, f, v) = R_B(s, f, v_B) = R_B(s, f, -v_A)$ and $\gamma_B(s, f, v)$ can be made dependent only on v_A ; in other words, it satisfies (2.7). It is then possible to define the operator \bar{T}_A^f as in (2.9). This operator maps \mathcal{W} to \mathcal{W} .

On the other hand, $(\bar{T}_A^f v_A)(s) = (T_A^{f, \gamma(s, f, v_A)} v_A)(s) = (T_A^{f, g} v_A)(s)$ for all $g \in R_B(s, f, v_A)$. But for every f, g , $T_A^{fg} v_A = -T_B^{fg} v_B$. And by definition (2.1), for all $g \in R_B(s, f, v_A)$, $h \in W_B$, $(T_B^{fg} v_B)(s) \geq (T_B^{fh} v_B)(s)$, which is equivalent to $(T_A^{fg} v_A)(s) \leq (T_A^{fh} v_A)(s)$. In other words,

$$R_B(s, f, v_A) = \arg \min_{g \in W_B} (T_A^{fg} v_A)(s)$$

so that

$$(\bar{T}_A^f v_A)(s) = \min_{g \in W_B} (T_A^{fg} v_A)(s).$$

As it was the case in (2.10), the minimization set in the right-hand side does not depend on v_A . The proof of Theorem 1 then applies mutatis mutandis, to conclude that operator \bar{T}_A is contractive on \mathcal{W} . It then admits a fixed point v_A^* in that set. Then the argument in the proof of Theorem 5 applies, and the FPE of the game is constructed from this fixed point. \square

3.4 Team Games

A result of [16] (Proposition 1) is that Team Games have an SSSE. Team Games (also known as Identical-Goal Games in [15]) are such that both players seek to maximize the same metric. This is a property of reward functions only. We slightly generalize the definition of these games and state a similar result for FPE.

Definition 5 (Team Game). The game is a Team Game if $\beta_A = \beta_B$ and there exists real constants μ and $\nu > 0$ such that: $r_B^{ab}(s) = \mu + \nu r_A^{ab}(s)$.

Theorem 8. If the game \mathcal{G} is a Team Game, then it admits an FPE.

Proof. We adapt the proof of Theorem 7. Consider a function v in the set $\mathcal{W} = \{v \in \mathcal{F}(\mathcal{S}) \times \mathcal{F}(\mathcal{S}) \mid v_B = \frac{\mu}{1-\beta} + \nu v_A\}$. Since v_B can be expressed as a function of v_A , then $SR_B(s, f, v)$ and $\gamma_B(s, f, v)$ can be made dependent only on v_A ; in other words, it satisfies (2.7). It is then possible to define the operator \bar{T}_A^f as in (2.9), mapping $\mathcal{F}(\mathcal{S})$ to $\mathcal{F}(\mathcal{S})$.

On the other hand, $(\bar{T}_A^f v_A)(s) = (T_A^{f, \gamma(s, f, v_A)} v_A)(s) = (T_A^{f, g} v_A)(s)$ for all $g \in R_B(s, f, v_A)$. A straightforward calculation concludes that for every f, g , $T_B^{fg} v_B = \frac{\mu}{1-\beta} + \nu T_A^{fg} v_A$. The operator T^{fg} then maps \mathcal{W} to \mathcal{W} and so does \bar{T}_A^f . And by definition (2.1), for all $g \in R_B(s, f, v_A)$, $h \in W_B$, $(T_B^{fg} v_B)(s) \geq (T_B^{fh} v_B)(s)$, which is equivalent to $(T_A^{fg} v_A)(s) \geq (T_A^{fh} v_A)(s)$ since $\nu > 0$. In other words,

$$R_B(s, f, v_A) = \arg \max_{g \in W_B} (T_A^{fg} v_A)(s)$$

so that

$$(\bar{T}_A^f v_A)(s) = \max_{g \in W_B} (T_A^{fg} v_A)(s) .$$

As it was the case in (2.10), the maximization set in the right-hand side does not depend on v_A . The proof concludes as in the proof of Theorem 7. \square

3.5 Acyclic Games

Acyclic games are such that state-to-state transitions do not lead back to a visited state, except for absorbing states. Acyclicity is a property of transition operators only.

We say that state s' is reachable from state s if there exist $k \in \mathbb{N}$, a sequence of states $s = s_0, s_1, \dots, s_k = s'$ and actions $a_0, \dots, a_{k-1}, b_0, \dots, b_{k-1}$ with $Q^{a_0 b_0}(s_1|s_0) \times Q^{a_1 b_1}(s_2|s_1) \times \dots \times Q^{a_{k-1} b_{k-1}}(s_k|s_{k-1}) > 0$.

Definition 6 (Acyclic Games). The game is an Acyclic Game if the state space \mathcal{S} admits the partition $\mathcal{S} = \mathcal{S}_\perp \cup \mathcal{S}_1$, with:

- for all $s \in \mathcal{S}_\perp, a \in \mathcal{A}_s, b \in \mathcal{B}_s, Q^{ab}(s|s) = 1$;
- for every pair $(s, s') \in \mathcal{S}_1 \times \mathcal{S}_1$, if s' is reachable from s , then s is not reachable from s' .

The following theorem is based on Theorem 4 and generalizes it for the FPE part.

Theorem 9. If the stochastic game \mathcal{G} is an Acyclic Game, then it admits an FPE.

Proof. The proof will proceed by successive reductions to static (or single-state) games. The game being acyclic, it is possible to perform a topological sort of the state space. There exists a partition $\mathcal{S} = \cup_{k=0}^K \mathcal{S}_k$ with $\mathcal{S}_0 = \mathcal{S}_\perp$ and for every $s \in \mathcal{S}_k, k > 0$, if s' is reachable from s then $s' \in \mathcal{S}_{k'}$ with $k' < k$. In a first step, we construct a candidate strategy (f^*, g^*) . Then we prove that this strategy solves the FPE problem.

For each $s_0 \in \mathcal{S}_0 = \mathcal{S}_\perp$, consider G_0 , the single-state game with $S = \{s_0\}$ and same strategies, rewards and discount factors. Theorem 4 applies to this game. It states that an FPE exists, resulting in a pair of strategies $(f_{s_0}^*, g_{s_0}^*) \in \mathbb{P}(\mathcal{A}_{s_0}) \times \mathbb{P}(\mathcal{B}_{s_0})$ and a value $V^*(s_0)$.

We now construct the strategies $(f_{s_k}^*, g_{s_k}^*)$ for $s_k \in \mathcal{S}_k$ with a recurrence on k . Assume this has been done up to $k-1$. Pick $s_k \in \mathcal{S}_k$. Because the game is acyclic, we have, for any $(f, g) \in W_A \times W_B$,

$$V_i^{fg}(s_k) = r_i^{fg}(s_k) + \beta_i \sum_{\ell=0}^{k-1} \sum_{s' \in \mathcal{S}_\ell} Q^{fg}(s'|s) V_i^*(s') . \quad (3.4)$$

Consider the static game (i.e. one-state game with null discount factors) with $S = \{s_k\}$ and rewards defined with this formula. Again, Theorem 4 applies to this game: an FPE exist, resulting in a pair of strategies $(f_{s_k}^*, g_{s_k}^*) \in \mathbb{P}(\mathcal{A}_{s_k}) \times \mathbb{P}(\mathcal{B}_{s_k})$. When $k = K$, we have defined this way a strategy (f_s^*, g_s^*) for each $s \in S$.

We now prove that this strategy is an FPE. We prove this with a recurrence. More precisely, we prove that for all k , property P_k holds, which says that for and all $s_k \in \mathcal{S}_k$:

$$\begin{aligned} g^*(s_k) &= \gamma_B(s_k, f^*, V^*) \\ f^*(s_k) &= R_A(s_k, V^*) \end{aligned}$$

where $V_i^* = V_i^{f^*g^*}$ is the unique fixed point of operator $T_i^{f^*g^*}$ for $i = A, B$. With Definition 3, the result will follow.

When $s_0 \in \mathcal{S}_0$, the local reaction set $LR_B(s_0, f, V^*)$ does not depend on V^* and $\{g_{s_0}^*\} \in LR_B(s_0, f^*, V^*)$, as in the proof of Theorem 4. In particular, it follows that $g^*(s_0) = \gamma_B(s_0, f^*, V^*)$ and $f^* = R_A(s_0, V^*)$. So P_0 holds.

Assume now that property P_ℓ holds for all $\ell < k$. Let $s_k \in \mathcal{S}_k$. Then

$$(T_B^{f^*g} V_B^*)(s_k) = r_i^{f^*g}(s_k) + \beta_i \sum_{\ell=0}^{k-1} \sum_{s' \in \mathcal{S}_\ell} Q^{f^*g}(s'|s) V_B^*(s'),$$

to be compared with (3.4). Then since $(f_{s_k}^*, g_{s_k}^*)$ solves (locally) the SSE for the subgame defined by (3.4), then $g_{s_k}^* = \gamma_B(s_k, f^*, V^*)$ and $f_{s_k}^* \in R_A(s_k, V^*)$ by construction. So property P_k holds. By recurrence, P_K holds and (f^*, g^*) is an FPE. \square

In contrast with Theorem 4, the existence of SSSE is not guaranteed for acyclic games. In Section 4.3 we study a game without an SSSE. This game is not acyclic, but it is possible to “approximate” it with an acyclic game which will have the same qualitative properties. On the other hand, if the transitions of a game are *deterministic*, in other words if the game is a multi-stage game, then it is MFS and it does have an SSSE according to Theorem 6.

4 Numerical Examples.

In this section, we present examples illustrating different situations that can occur by comparing the solutions returned by the different algorithms presented in Section 2.

In the example of Section 4.2, VI converges to some FPE. Also we show that the FPE is an SSSE and the solution returned by (MP). The model involved does not satisfy the sufficient conditions for existence and convergence identified in Sections 2.3 and 3. In Section 4.3, we describe an example where, depending on the discount factors β_A, β_B : either an SSSE exists and does not coincide with the FPE, or an SSSE does not exist, or an FPE does not exist. Finally, in Section 4.4, we describe an example where an FPE is shown to exist, but VI does not necessarily converge to it.

4.1 Experimental setup

The numerical experiments involving Value Iteration and Policy Iteration were performed using Python 3.6, on a machine running under MacOS, a processor of 2,6 GHz Intel Core i5, and memory of 8 GB 1600 MHz DDR3. Operator T is implemented with Cplex 12.8.

To solve the linear systems involved in T^{fg} we use the Python library Numpy. In order to solve (MP) we use the KNITRO 12.0 solver combined with AMPL.

The data of the experiments is presented tables with the convention of the following figure. It shows the relevant parameters when the system is in state s , Player A performs action a and Player B performs action b :

	b	
a	$(Q^{ab}(s_1 s), Q^{ab}(s_2 s))$	$(r_A^{ab}(s), r_B^{ab}(s))$

4.2 Example 1: FPE and SSSE coincide and VI converges

This first example shows the convergence of Value Iteration in a case where an FPE exists. Consider $\beta_A = \beta_B = \frac{9}{10}$ and the data in Table 1.

	b_1		b_2	
a_1	$(\frac{1}{2}, \frac{1}{2})$	$(10, -10)$	$(0, 1)$	$(-5, 6)$
	$(\frac{1}{4}, \frac{3}{4})$	$(-8, 4)$	$(1, 0)$	$(6, -4)$
State s_1				
	b_1		b_2	
a_1	$(\frac{1}{2}, \frac{1}{2})$	$(7, -5)$	$(0, 1)$	$(-1, 6)$
	$(\frac{1}{4}, \frac{3}{4})$	$(-3, 10)$	$(1, 0)$	$(2, -10)$
State s_2				

Table 1: Transition matrix and payoffs for each player.

This example does not satisfy any of the sufficient conditions listed in Sections 2 and 3. However, both Value Iteration and Policy Iteration converge to the FPE, and furthermore the FPE, the SSSE and the optimal solution returned by (MP) coincide. The application of Value Iteration, starting with the null function, results in the evolution displayed in Figure 1. Given that Value Iteration converges, an FPE exists. The policies and values are given in Table 2. Detailed algebraic manipulations of the value functions given the data in this example allows us to show that the SSSE satisfies the values in Table 2. Then, we show that this solution is a fixed point for operator T . These details are provided in Appendix C.1.

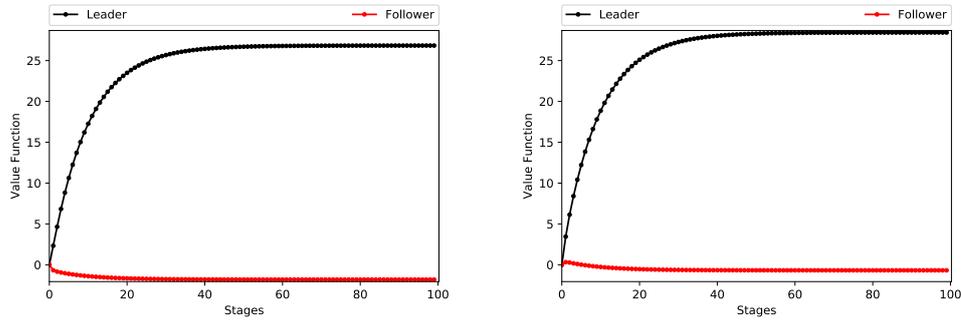


Figure 1: Value Iteration applied to Example 1

	s_1	s_2
Play of A	(0.3467, 0.6533)	(0.6434, 0.3566)
Play of B	b_1	b_2
v_A	26.841	28.437
v_B	-1.807	-0.679

Table 2: Policies and values of the SSSE in Example 1

4.3 Example 2: FPE and SSSE are different

In this section we will study the stochastic game given by the data in Table 3. This game has two states and two actions per state. Figure 2 shows a diagram of the possible transition between states. State s_1 is absorbing for any combination of actions.

	b_1	b_2		b_1	b_2
a_1	(1, 0) (-1, -2)	(1, 0) (-2, 1)	State s_1	a_1	(1, 0) (-1, -2) (-2, -2)
a_2	(1, 0) (0, 0)	(1, 0) (2, 0)		a_2	(0, 1) (0, 1) (1, 0) (1, 1)

Table 3: Transition matrix and payoffs for each player in Example 2

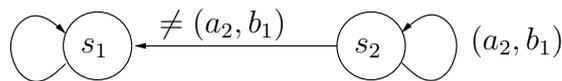


Figure 2: Transition structure of Example 2

Whenever $\beta_B > 0$, depending on the values of β_A the existence and non-existence

for both equilibrium concepts, FPE and SSSE, changes. For $\beta_A < \frac{1}{5}$ the existence of an FPE is guaranteed, but no SSSE exist. This comes from the fact that for these values of β_A , $R_A(s_1) = \{a_1\} \times \{(f_2, 1 - f_2) : f_2 \in [0, 1]\}$ and $R_A(s_2) = \{a_2\} \times \{a_2\}$. Clearly, $R_A(s_1) \cap R_A(s_2) = \emptyset$ and therefore there is no SSSE. When the game starts in the absorbing state s_1 , the optimal policy for the leader is to play and announce the (static) SSE in s_1 and an arbitrary strategy in s_2 . When the game starts in s_2 , the leader has the incentive to announce a sub-optimal strategy in s_1 in order to remain in state s_2 and increase his expected reward. This could be done only for low values of β_A .

In the interval $\beta_A \in [\frac{1}{5}, \frac{1}{3}]$, both SSSE and FPE exists, but the strategies and values are different. On the other hand, whenever $\beta_A > \frac{1}{3}$, there is no FPE. This information is summarized in Table 4. The details of this analysis are provided in Appendix C.2. We explain in particular, in Section C.2.3.2.1, the complex dynamics of the operator when $\beta_A > \frac{1}{3}$. The dynamics of Value Iteration is numerically illustrated in Figure 3. In the top two graphs, $\beta_A > 1/3$ and the iterations for state s_2 exhibit a cycle of order 3. In the bottom two graphs, $\beta_A < 1/3$ and Value Iteration converges for both states to the FPE identified in Table 4.

β_A	SSSE				FPE			
	$v_A(s_1)$	$v_A(s_2)$	$v_B(s_1)$	$v_B(s_2)$	$v_A(s_1)$	$v_A(s_2)$	$v_B(s_1)$	$v_B(s_2)$
$[0, \frac{1}{5})$	No	No	No	No	$\frac{2}{1-\beta_A}$	0	0	$\frac{1}{1-\beta_B}$
$[\frac{1}{5}, \frac{1}{3}]$	$\frac{2}{1-\beta_A}$	$\frac{2\beta_A}{1-\beta_A}$	0	0	$\frac{2}{1-\beta_A}$	0	0	$\frac{1}{1-\beta_B}$
$(\frac{1}{3}, 1)$	$\frac{2}{1-\beta_A}$	$\frac{2\beta_A}{1-\beta_A}$	0	0	No	No	No	No

Table 4: Existence of FPE and SSSE for different values of β_A and $\beta_B > 0$. Even when both exist they may not coincide.

Policy iteration and (MP). Finally, we test the Mathematical Programming (MP) formulation and Policy Iteration (PI) algorithm, for the different values of the parameter α in (MP) and for values of β_A in $\{0, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}\}$. Table 5 summarizes the results obtained.

In this experiment whenever the SSSE exists ($\beta_A \in \{\frac{1}{4}, \frac{1}{2}\}$), (MP) computes it correctly and Policy Iteration converges to the FPE. When no SSSE exists ($\beta_A \in \{0, \frac{1}{8}\}$), (MP) returns a value that is influenced by the vector of weights α_s . In these cases Policy iteration converges to an FPE that is different from the solution found by (MP).

4.4 Example 3: FPE exists but VI does not converge to it

We now develop an example where an FPE does exist, but Value Iteration does not necessarily converge to it. The data of this example is listed in Table 6.

Consider $\beta_A = \beta_B = \frac{1}{2}$. We claim that the pair of strategies (f^*, g^*) and value functions (v_A^*, v_B^*) in Table 7 constitute *both* an SSSE and an FPE. We provide in Appendix C.3 justifications for this claim.

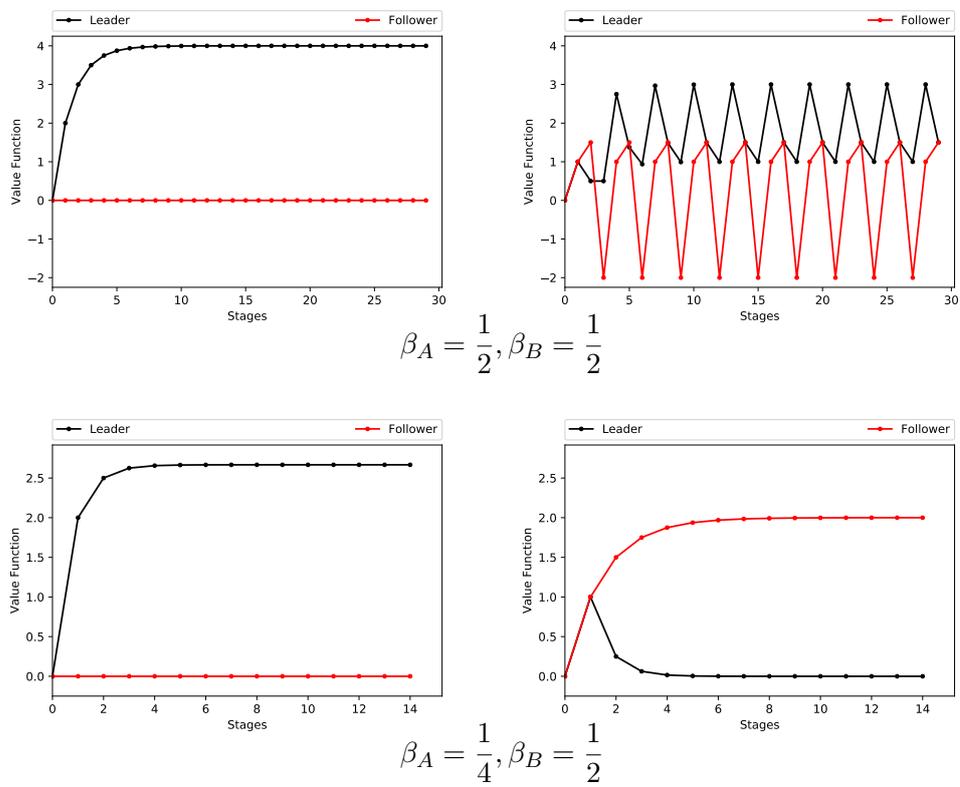


Figure 3: Value Iteration applied to Example 2: state s_1 (left) and s_2 (right)

β_A		(MP)						(PI)	
		$\alpha_{s_1} = 100$	$\alpha_{s_2} = 1$	$\alpha_{s_1} = 1$	$\alpha_{s_2} = 100$	$\alpha_{s_1} = 1$	$\alpha_{s_2} = 1$	s_1	s_2
0	v_A	2	~ 0	-2	1	2	~ 0	2	0
	v_B	~ 0	~ 0	2	2	~ 0	~ 0	0	2
	f	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(1,0)	(0,1)	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	(0,1)
	g	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(1,0)
$\frac{1}{8}$	v_A	16/7	2/7	-16/7	5/7	16/7	2/7	16/7	0
	v_B	~ 0	~ 0	2	2	~ 0	~ 0	0	2
	f	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(1,0)	(0,1)	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	(0,1)
	g	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(1,0)
$\frac{1}{4}$	v_A	8/3	2/3	8/3	2/3	8/3	2/3	8/3	0
	v_B	~ 0	~ 0	~ 0	~ 0	~ 0	~ 0	0	2
	f	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	(0,1)
	g	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(1,0)
$\frac{1}{2}$	v_A	4	2	4	2	4	2	-	-
	v_B	~ 0	~ 0	~ 0	~ 0	~ 0	~ 0	-	-
	f	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	(0,1)	$(\frac{1}{3}, \frac{2}{3})$	-	-
	g	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	(0,1)	-	-

Table 5: Results for (MP) and (PI) with different values of α and β_A with $\beta_B = 0.5$ fixed

	b_1	b_2		b_1	b_2
a_1	(1, 0)	(0, 1)		(0, 1)	(1, 0)
	(1, -1)	(0, 1)		(-1, 0)	(0, 1)
a_2	(0, 1)	(0, 1)		(1, 0)	(0, 1)
	(-1, 1)	(-1, -1)		(0, 1)	(1, -1)
	State s_1			State s_2	

Table 6: Transition matrix and payoffs for each player in Example 3

	s_1	s_2
Play of A	(1, 0)	$(5 - \sqrt{19}, -4 + \sqrt{19})$
Play of B	b_2	b_2
v_A	$\frac{1}{5}(-3 + \sqrt{19})$	$\frac{1}{5}(-6 + 2\sqrt{19})$
v_B	$\frac{1}{5}(16 - 2\sqrt{19})$	$\frac{1}{5}(22 - 4\sqrt{19})$

Table 7: Values and Policies forming an SSSE and FPE.

When applying Value Iteration with the null function as a starting point, we get however the evolution in Figure 4. Values obtained with Policy Iteration have a similar behavior. Finally, (MP) returns as the optimal solution the SSSE (and FPE).

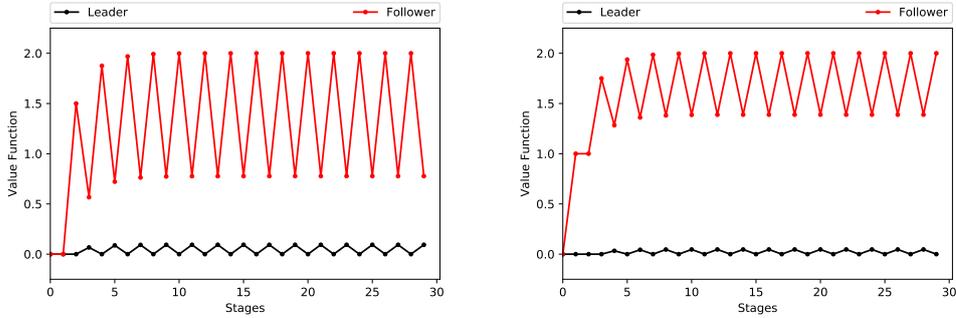


Figure 4: Value Iteration applied to Example 3: state s_1 (left) and s_2 (right)

5 Application: Surveillance in a graph.

In this section we present an example of a stochastic game that models the interaction between a security patrol and an attacker. In this game a defender has to patrol (or “cover”) a set of locations and an attacker wants to perform an attack in one of these locations, both maximizing their expected rewards.

The state of the game, known to both players, is defined by the locations of the defender and attacker and whether an attack occurred or not. Given this information, the leader decides on a patrolling strategy (a location to patrol) which is also known to the attacker when deciding two things: where to move and if he performs an attack or not. Once the attack is performed the game ends in one of two terminal states. This game model situations where the defender has the information where the possible attacker it is located, but they cannot perform any action if an attack is not performed. In real situations, it can be the case of demonstrations or high-risk football matches.

The rewards mainly depend on the place where the attack is performed and whether the location is being covered or not. The effectiveness of the player’s movements is influenced by random factors: when they decide to move to some location this move may fail due to external factors, in which case they remain in their current location.

5.1 Game description

We introduce now the elements of the model and the notation. A synthesis of this notation is presented in Appendix A. Formally, we consider a set of locations to patrol/targets $\mathcal{L} = \{l_1, l_2, \dots, l_n\}$. There are some connections allowed between locations represented by edges (denoted by \mathcal{E}), so the board of the game is actually a graph $(\mathcal{L}, \mathcal{E})$. Player A is the defender, Player B is the attacker. The state space is $\mathcal{S} = \mathcal{L} \times \mathcal{L} \times \{0, 1\} \cup \{\perp_0, \perp_1\}$.

A typical state $s = (\ell_A, (\ell_B, \alpha)) \in \mathcal{S}$ represents the defender's location ($\ell_A \in \mathcal{L}$) and the attacker's location ($\ell_B \in \mathcal{L}$). The binary parameter α takes the value 1 if the attacker is committing an attack or 0 if he is not, thereby being unnoticed in that period. There are also two special fictitious states \perp_0, \perp_1 representing the state of the game once the attack was performed. \perp_1 represents the case where the attack is successful and \perp_0 when the attacker is caught. These two are absorbing states in our game.

The action space $\mathcal{A}_s \subset \mathcal{L}$ for the leader represents all the possible location that he can achieve from its current position (given by the state s). For the follower, $\mathcal{B}_s \subset \mathcal{L} \times \{0, 1\}$ represents all the possible locations that the attacker can achieve in the state s with and the decision whether to attack or stay unnoticed. We use the notation $\ell \in \mathcal{L}$ to represent the action of “move to ℓ ” and $\alpha \in \{0, 1\}$ to represent the action of “attacking” or “stay unnoticed” respectively. In states $s \in \{\perp_0, \perp_1\}$, actions are irrelevant: we can pick $\mathcal{A}_s = \mathcal{L}$ and $\mathcal{B}_s = \mathcal{L} \times \{0, 1\}$ for such states by convention.

The probability of transitions between states is constructed using the function $q_i^{\ell'}(\ell''|\ell)$ which denotes the probability that player $i \in \{A, B\}$ reaches location ℓ'' from ℓ having decided to move to ℓ' . We assume that these probabilities are independent between players. In particular, if there is no failure in Player i 's movements, then $q_i^{\ell'}(\ell''|\ell) = 1$ if $\ell'' = \ell'$, and 0 otherwise. The transition probabilities $Q^{ab}(z|s)$ are then defined from these quantities by expressions (5.1)–(5.3) as follows:

$$Q^{ab}(z|\perp_i) = \begin{cases} 1 & z = \perp_i, \quad i \in \{0, 1\}, (a, b) \in \mathcal{A}_{\perp_i} \times \mathcal{B}_{\perp_i} \\ 0 & \text{otherwise.} \end{cases} \quad (5.1)$$

$$Q^{\ell'_A, (\ell'_B, \alpha)}(\ell''_A, (\ell''_B, \alpha)|\ell_A, (\ell_B, 0)) = q_A^{\ell'_A}(\ell''_A|\ell_A) q_B^{\ell'_B}(\ell''_B|\ell_B) \quad \begin{array}{l} \alpha \in \{0, 1\}, \ell'_A \in \mathcal{A}_{\ell_A}, \\ (\ell'_B, \alpha) \in \mathcal{B}_{\ell_B}, \end{array} \quad (5.2)$$

$$Q^{\ell'_A, (\ell'_B, \alpha)}(z|\ell_A, (\ell_B, 1)) = \begin{cases} 1 & z = \perp_0 \text{ and } \ell_A = \ell_B \\ 1 & z = \perp_1 \text{ and } \ell_A \neq \ell_B \\ 0 & \text{otherwise.} \end{cases} \quad (5.3)$$

Rewards result from the interaction, or lack thereof, between the defender and the attacker. Following the notation used in security games, we denote $U_A^u(\ell) < 0$ and $U_A^c(\ell) > 0$ the penalty and the benefit for the defender if an attack is performed in ℓ , which depends only if the target is uncovered (superscript u) or covered (superscript c). Similarly, we define $U_B^u(\ell) > 0$ and $U_B^c(\ell) < 0$ the reward (and penalty respectively) of the attacker if the location attacked is uncovered or not. Instant rewards r_A and r_B are defined as the expected values of the rewards $R_i = R_i^{ab}(z|s)$, $i \in \{A, B\}$, of the dynamics between players, which depends on the current state of the system s , the actions (a, b) performed by the players and the future state of the system z . This technique is fairly standard, as it is shown in [13, Ch. 2, pp.20]. The expressions for $R_i^{ab}(z|s)$, $i = A, B$ are

listed in (5.4), in which s is any state and $z = (\ell_A, (\ell_B, \alpha))$.

$$\begin{array}{rcc}
 R_A^{ab}(z|s) & R_B^{ab}(z|s) & \\
 \hline
 U_A^u(\ell_B) & U_B^u(\ell_B) & \ell_A \neq \ell_B \text{ and } \alpha = 1 \\
 U_A^c(\ell_B) & U_B^c(\ell_B) & \ell_A = \ell_B \text{ and } \alpha = 1 \\
 P_A(\ell_A) & P_B(\ell_B) & \alpha = 0 \\
 P_A(\perp_0) & P_B(\perp_0) & z = \perp_0 \\
 P_A(\perp_1) & P_B(\perp_1) & z = \perp_1 \\
 \hline
 \end{array} \tag{5.4}$$

The first two lines represent the payoffs when the attack is performed, $z = (\ell_A, (\ell_B, 1))$. In the third line, $P_B(\ell_B) < 0$ represents the opportunity cost and risk for the attacker of being in location ℓ_B and not perform an attack. In the two last lines P_A and P_B represents the residual value of being in an absorbing state. We assume $P_A(\perp_0) > 0$ and $P_A(\perp_1) < 0$ and the opposite for the defender. From the definition of R_A and R_B we obtain instants rewards r_A and r_B as follows:

$$r_i^{ab}(s) = \sum_{z \in \mathcal{S}} Q^{ab}(z|s) R_i^{ab}(z|s) \quad i \in \{A, B\}.$$

The dynamics of the game is summarized as follow: First, at the start of any epoch, the system is in a state formed by the location the both players and the behaviour of the attacker. Then, the defender knowing the state of the game chooses a strategy f (probably mixed) over the locations reachable from his current location. The attacker observes the strategy and chooses where to move and whether to attack or not. We denote this action as g . Note, that if the attacker decide to attack, the success or failure of his strategy will be revealed in the next state. The system evolves to the following state influenced by f , g , and Q . Both players receive their payoffs.

5.2 Computational study

Here we evaluate the solution algorithms presented in terms of solution time and quality of the solution obtained. We begin by describing the instances of the graph surveillance problem constructed for this computational study. We compare the solution times of Value Iteration and Policy Iteration on every instance considered. Solving the (MP) formulation was not a competitive solution approach for these problems. Even for the smallest instance considered, solving the (MP) formulation using a state-of-the-art non-linear optimization package did not return a solution after more than 5 hours of computational time, which far exceeds the solution times observed for Value Iteration and Policy Iteration over all instances. We therefore do not present computational results of solving the (MP) formulation. The experimental setup used is the one described in Section 4. To evaluate the quality of the FPE solutions obtained by these algorithms we compare them to heuristic policies, both in the myopic and non-myopic follower case.

We generate instances of different structure and size by considering different graphs to patrol on: paths, cycles, T-shaped graphs and complete graphs. We limit the size of \mathcal{A}_s and \mathcal{B}_s by limiting the distance that each player can travel from one time step to the

next. To do so, we introduce the parameter k as the maximum geodesic distance that each player can travel through one time step.

Functions $q_A^{\ell'}(\ell''|\ell)$ are a function of the nodes that are in the shortest path between ℓ and ℓ' . We denote this set of nodes as $SP(\ell, \ell') = \{\ell, \ell_{i_2}, \dots, \ell_{i_{k-1}}, \ell'\}$. Probabilities q_A are defined as follows: if $|SP(\ell, \ell')| = 1$ and $SP(\ell, \ell) = \{\ell\}$, then $q_A^{\ell'}(\ell|\ell) = 1$, if $|SP(\ell, \ell')| \geq 2$ then

$$q_A^{\ell'}(\ell''|\ell) = \begin{cases} 1 - \epsilon & \ell' = \ell'' \\ \frac{\epsilon}{|SP(\ell, \ell')| - 1} & \ell'' \in SP(\ell, \ell') \setminus \{\ell'\} \\ 0 & \text{otherwise.} \end{cases} \quad (5.5)$$

In our experiments we set the probability of failing to $\epsilon = 0.25$. We assume $q_B^{\ell'}(\ell''|\ell) = \mathbf{1}_{\ell'=\ell''}$ are deterministic: the attacker always succeeds with its intended move.

The payoff functions are defined in Table 8. The values of each parameter depend on the degree of the node, representing the fact that nodes with greater degree are more important in order to keep the connectivity of the graph.

Parameter	Value	Parameter	Value
$U_A^u(\ell)$	$-10deg(\ell)$	$U_A^c(\ell)$	$10deg(\ell)$
$U_B^u(\ell)$	$2^{deg(\ell)}$	$U_B^c(\ell)$	$-2^{deg(\ell)}$
$P_A(\ell_A)$	0	$P_B(\ell_B)$	$-deg(\ell_B)$
$P_A(\perp_0)$	1	$P_A(\perp_1)$	-1
$P_B(\perp_0)$	-1	$P_B(\perp_1)$	1

Table 8: Payoff functions description

We test our models in instances with $n \in \{5, 10\}$, $\beta_B = \{0, 0.9\}$ and $k \in \{2, 3\}$ for each type of graph. We show how the instances increase with n and k in Table 9, which gives the size of the state set and the average size of the sets of actions for every instance considered. The stopping criterion is set to $\epsilon = 10^{-3}$ for both Value Iteration and Policy Iteration.

The experimental setup is the one described in Section 4.1. Figure 5 shows the solution times in a performance profile in logarithmic scale comparing Value Iteration and Policy Iteration over all the instances considered. Policy iteration has faster solution times over the instances tested. We present the solution time results separated for each instance type in Figure 6. The results show that graph structure does not influence significantly the solution times of Value Iteration and Policy Iteration.

We now aim to evaluate the quality of the Stackelberg equilibrium solution for this discounted game. However, since we do not have a method to compute the SSSE in general, we use the FPE that can be computed using Value Iteration as a proxy. Accordingly, we compute the values of equilibrium $v^* = (v_A^*, v_B^*)$, with the respective equilibrium policies f^* and g^* using the Value Iteration algorithm.

We compare the FPE solution to static policies obtained by ignoring the dynamic nature of the game. We refer to these heuristic policies as Myopic policies. To determine

n	k	type	$ \mathcal{S} $	$ \mathcal{A} $	$ \mathcal{B} $
5	2	Cycle	52	4.8	9.7
		Line	52	3.7	7.3
		T	52	4.1	8.1
		Complete	52	4.8	9.7
	3	Cycle	52	4.8	9.7
		Line	52	4.5	8.9
		T	52	4.8	9.7
		Complete	52	4.8	9.7

n	k	type	$ \mathcal{S} $	$ \mathcal{A} $	$ \mathcal{B} $
10	2	Cycle	202	5.0	9.9
		Line	202	4.4	9.9
		T	202	4.6	8.7
		Complete	202	9.9	9.1
	3	Cycle	202	6.9	13.9
		Line	202	5.8	11.5
		T	202	6.3	12.7
		Complete	202	9.9	19.8

Table 9: Size of the instances (number of states and average number of actions for Player A and B) when changing the graph structure, number of nodes n and displacement limit k

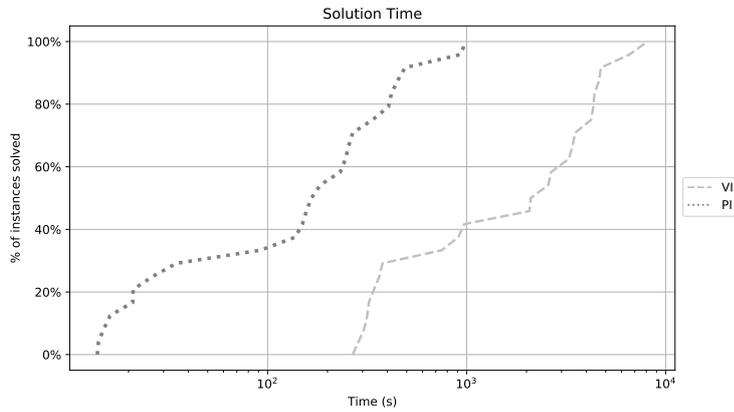


Figure 5: Solution time performance profile.

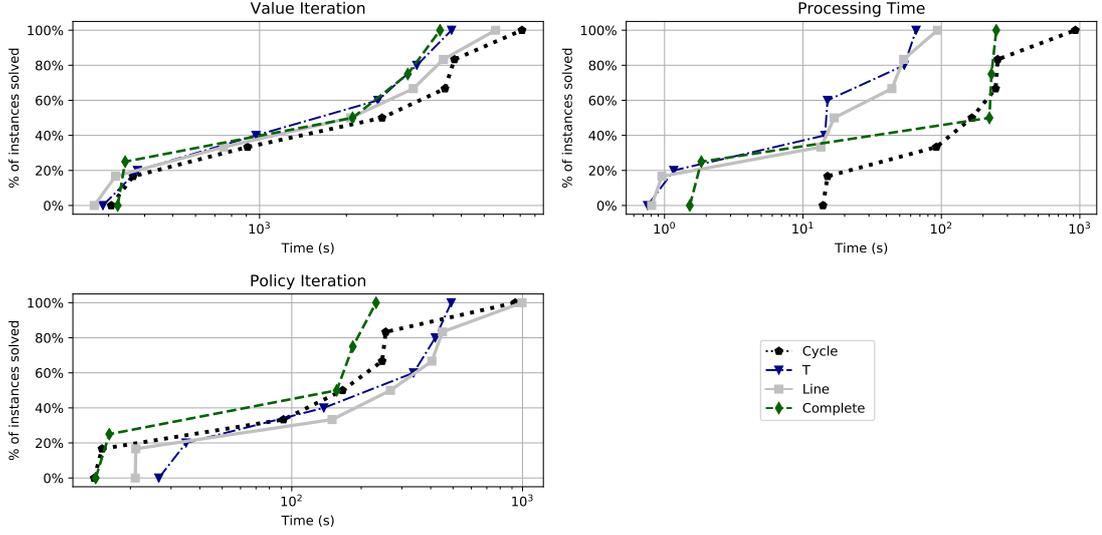


Figure 6: Solution time performance profile for each instance type.

the Myopic policy, for each state we compute the Strong Stackelberg policies, f^M, g^M of the static game, that is with $\beta_A = \beta_B = 0$. Finally, we evaluate this policy in the dynamic setting: we obtain the value $V^{f^M g^M} = v^M = (v_A^M, v_B^M)$ as the fixed point of operator $T^{f^M g^M}$ (see (2.6)) with real values of β_A and β_B .

Finally, in order to compare the policies obtained for both methods we compare the average value for applying each policy, denoted respectively as $\bar{v}_A^*, \bar{v}_B^*, \bar{v}_A^M$ and \bar{v}_B^M , where $\bar{v} = |\mathcal{S}|^{-1} \sum_{s \in \mathcal{S}} v(s)$.

Table 10 shows the comparison of the values for the different types of graph structures mentioned before, and with the parameters $n = 10, k = 2, \beta_A = 0.9$ and $\beta_B = 0$. Recall that in this case, Value Iteration converges to an FPE (and SSSE) because the game is MFS. For the complete graph, the myopic strategy generates in average the same reward as the equilibrium strategy. In the other cases, the SSSE strategy outperforms the myopic-heuristic policy.

type	\bar{v}_A^*	\bar{v}_A^M	\bar{v}_B^*	\bar{v}_B^M
Cycle	9.957	8.376	1.485	2.079
Path	9.070	6.686	1.109	1.703
T	10.623	8.129	0.703	2.218
Complete	89.595	89.595	129774.653	129774.653

Table 10: Evaluation of the solution concept with $\beta_B = 0$.

We repeat the same evaluation, but now with $\beta_B = 0.9$. Note that in this case there are no guarantees that Value Iteration will converge, or that the FPE policy is an SSSE. Value Iteration found an FPE (i.e. converged) in all the instances. Table 11 shows

the average values obtained applying the policies provided by the FPE and the myopic heuristic.

type	\bar{v}_A^*	\bar{v}_A^M	\bar{v}_B^*	\bar{v}_B^M
Cycle	-17.667	6.767	6.261	2.171
Path	-14.102	4.938	6.617	1.506
T	-12.955	6.143	6.739	2.249
Complete	89.595	89.604	129773.771	129773.762

Table 11: Evaluation of the solution concept with $\beta_B = 0.9$.

Note that in relative terms, both the FPE policy and the Myopic policy return a lower average value for Player A. This can be attributed to the follower's change of behavior. More significantly, the FPE solution found obtains lower values for the leader than the Myopic strategy, significantly smaller except for the Complete graph, which is smaller but comparable. In this case, with $\beta_B = 0.9$, the game does not satisfy MFS and the FPE solution obtained by the Value Iteration may not correspond to the SSSE, which may lead to inefficient solutions.

6 Conclusions and Further Work

In this paper, we have demonstrated the relevance of the concept of Strong Stationary Stackelberg Equilibria, SSSE, and the related operator-based algorithms, for the computation of policies in the context of two-player discounted stochastic games.

For this, we first defined a suitable operator acting on the set of value functions for both players. We introduced the concept of Fixed-Point Equilibrium, FPE, as the fixed points of this operator. We then investigated the relationship between SSSE and FPE. We show that neither need to exist in general, and that when they do, they do not necessarily coincide. We also show that the solution based on Mathematical Programming suggested in the literature, does not necessarily compute a correct answer. We have nevertheless identified several classes of games where SSSE and FPE do exist and do coincide. Among these is the class of games with Myopic Follower Strategies, MFS, which include games with myopic followers and games with leader controlled transitions.

We consider an application in a security domain, in which a moving defender protects locations on a graph that can be attacked by a moving attacker. We give a formulation of this problem as a Stackelberg equilibrium in a discounted stochastic game. The Value Iteration and Policy Iterations algorithms are able to compute efficiently the FPE for the instances considered of this security application. The instances considered were too difficult to solve if using the Mathematical Programming formulation of the problem. In the case of myopic follower, in which the FPE corresponds to the SSSE, we observe that the solution obtained is efficient and outperforms heuristic policies. However, in examples without MFS, we see that the FPE solution computed is worse than heuristic policies.

Future research will aim at identifying more general sufficient conditions for the two concepts, SSSE and FPE, to coincide. The problem of finding general methodologies to detect the existence of SSSE is still open. It is also important to determine algorithms to find the equilibrium in games which possess SSSEs but do not satisfy the MFS condition. Finally, It will also be interesting to determine whether the use of Value Iteration or Policy iteration, when they do not converge, can nevertheless produce nonstationary strategies with a good performance.

Acknowledgments.

Victor Bucarey has been partially supported by the Fonds de la Recherche Scientifique -FNRS under Grant(s) no PDR T0098.18.

A Notation Summary

We list in the following tables the principal notation used in the paper. Table 12 groups the notation related to the general model. Table 13 is relative to the application of Section 5.

Stochastic Games	
\mathcal{G}	Stochastic game.
\mathcal{S}	Set of states.
$\mathcal{A}_s, \mathcal{B}_s$	Set of actions for the Leader and the Follower respectively available in state s .
$r_i^{ab}(s)$	Immediate reward for Player i when actions a and b are performed in state s .
$Q^{ab}(z s)$	Probability of reaching state z from state s when actions a and b are performed.
β_i	Discount factor for Player i .
W_i	Set of feedback policies for Player i .
f, g	Policies for the Leader (Player A) and the Follower (Player B) respectively.
$V_i^{fg}(s)$	Expected discounted sum of all rewards for Player i , when policies f and g is applied and the starting state is s .
R_B, SR_B	Set of optimal responses (and strong respectively) for the Follower.
γ_B	Best response which is performed for the Follower.
$R_A(s)$	Best strategies for the Leader starting in state s .
v_i	Value function for Player i .
$\mathcal{F}(\mathcal{S})$	Set of value functions $\mathcal{S} \rightarrow \mathbb{R}$.
T_i^{fg}	One-step operator for any fixed pair of policies f and g .
T	Dynamic programming operator.
$\bar{\gamma}_B$	Best response of the Follower when it does not depend on v_B .
\bar{T}_A^f, \bar{T}	One-step operator and Dynamic programming operator when T does not depend on v_B .
M_i	Upper bound in value functions for Player i .

Table 12: Table of general notation

Application: Surveillance in a graph	
\mathcal{L}, \mathcal{E}	Set of locations and connections between locations.
\perp_0, \perp_1	Absorbing states of the game.
α	Binary decision of attack or not.
$U_i^u(\ell), U_i^c(\ell)$	Reward or penalty for player i when an attack is performed in location ℓ and it is being protected (c) or unprotected (u).
P_i	Rewards and costs for player i when an attack is not performed.
$q_i^{\ell'}(\ell'' \ell)$	Probability for player i of reaching location ℓ'' from ℓ given that he decides to move to ℓ' .
n	Number of locations in the graph.
k	Maximum geodesic distance which the defender can move in each step time.
\bar{v}^*	Average value function when a the policy obtained via Value Iteration is performed.
\bar{v}^M	Average value function when a myopic policy is performed.

Table 13: Table of notation for the surveillance application

B Algorithmic Proof for the Static Case.

We provide an algorithmic proof for the existence of SSSE in Theorem 4. This proof is also useful because it allows to compute a SSE in the static case, which is an important step in Algorithms 1, 2 and 4.

Lemma 6. If the game \mathcal{G} is static (i.e. has only one state), then it has a SSSE.

Proof. Proof. Since the state space has only one state, we omit the reference to states in the notation. For each action $b \in \mathcal{B}$, define the following problem LP(b):

$$\text{LP}(b) \quad \max \sum_{a \in \mathcal{A}} r_A^{ab} f(a) \quad (\text{B.1})$$

$$\text{s.t.} \quad \sum_{a \in \mathcal{A}} r_B^{ab} f(a) \geq \sum_{a \in \mathcal{A}} r_B^{ab'} f(a) \quad \forall b' \in \mathcal{B} \quad (\text{B.2})$$

$$\sum_{a \in \mathcal{A}} f(a) = 1 \quad (\text{B.3})$$

$$f(a) \geq 0 \quad \forall a \in \mathcal{A}. \quad (\text{B.4})$$

Consider Algorithm 5. Each of these LPs is bounded. Therefore, they are either unfeasible or they have a finite optimal solution. Furthermore, at least one of them is feasible. In that case, the optimal value represents the maximum value achieved by incentivising to the follower to play b . By picking the maximum of these values, the optimal f of that LP with the action b is an SSE for the game. \square

Algorithm 5 Multiple LPs Algorithm

- 1: Initialize with $n = 0$, $v_A^0(s) = 0$ for every $s \in \mathcal{S}$
 - 2: **for** $b \in \mathcal{B}$ **do**
 - 3: Solve LP(b)
 - 4: **if** LP(b) is bounded **then**
 - 5: Retrieve $f(b)$ and $v^*(b)$, respectively the leader's optimal policy given that the follower plays b , and its value
 - 6: **end if**
 - 7: **end for**
 - 8: Set the policy g^* , in which the follower plays $b^* = \arg \max v^*(b)$
 - 9: **return** SSE $(f^*(b^*), g^*)$ for the static game
-

Given that Algorithm 5 solves $|\mathcal{B}|$ LPs with $|\mathcal{A}|$ variables and $|\mathcal{B}|$ constraints, it finds a SSSE in polynomial time.

C Analysis of the examples

In this section, we provide details and justifications for the examples reported in Section 4. In doing so, we introduce some elements useful in the analysis of the operators involved in Value Iteration algorithms.

C.1 Analysis of Example 1

In this section, we provide details on the analysis of Example 1 introduced in Section 4.2. We compute the Stationary Strong Stackelberg Equilibria (SSSE) and demonstrate their existence. We compute the one-step operator and show it has also a fixed point (FPE) which coincides with this SSSE.

C.1.1 Data

The data of this example is given in Table 14 (also Table 1). We will also specifically consider $\beta_A = \beta_B = 9/10$.

		b_1	b_2			b_1	b_2
a_1	$(\frac{1}{2}, \frac{1}{2})$	$(10, -10)$	$(0, 1)$	a_1	$(\frac{1}{2}, \frac{1}{2})$	$(7, -5)$	$(0, 1)$
a_2	$(\frac{1}{4}, \frac{3}{4})$	$(-8, 4)$	$(1, 0)$	a_2	$(\frac{1}{4}, \frac{3}{4})$	$(-3, 10)$	$(1, 0)$
		State s_1				State s_2	

Table 14: Transition matrix and payoffs for each player in Example 1

C.1.2 Computation of the SSSE

C.1.2.1 Values of stationary strategies. Given a stationary strategy (f, g) , let us compute the value V_i^{fg} of this strategy for each player $i = A, B$. The notation is $f_i = f(s_i, a_i)$ and $g_i = g(s_i, b_i)$. We have first, observing that transition probabilities do not depend on the state,

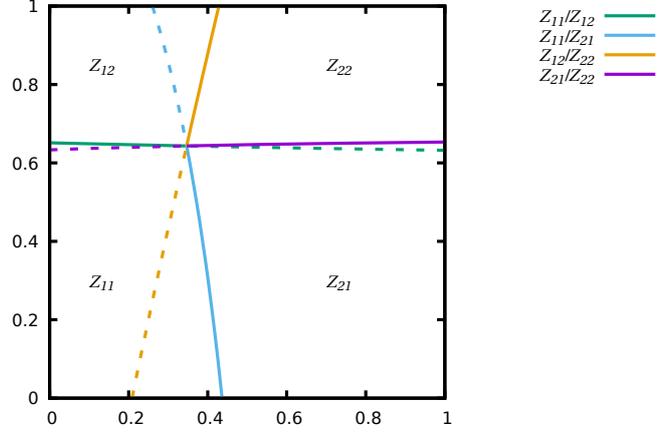
$$Q^{fg} = \begin{pmatrix} \frac{5}{4}f_1g_1 - f_1 - \frac{3}{4}g_1 + 1 & -\frac{5}{4}f_1g_1 + f_1 + \frac{3}{4}g_1 \\ \frac{5}{4}f_2g_2 - f_2 - \frac{3}{4}g_2 + 1 & -\frac{5}{4}f_2g_2 + f_2 + \frac{3}{4}g_2 \end{pmatrix}.$$

Also,

$$r_{AB}^{fg} = \begin{bmatrix} (29g_1 - 11)f_1 - 14g_1 + 6 & (-24g_1 + 10)f_1 + 8g_1 - 4 \\ (13g_2 - 3)f_2 - 5g_2 + 2 & (-31g_2 + 16)f_2 + 20g_2 - 10 \end{bmatrix}.$$

From these elements, one derives the values for Player B in each state, when she plays the four possible strategies: with $V_B^{f,xy}(s)$ the value at state s when x is played in state s_1 and y played in state s_2 :

$$V_B^{f,b_1b_1}(s_1) = 10 \frac{9f_1f_2 - 272f_1 - 369f_2 + 322}{-9f_1 + 9f_2 + 40} \quad V_B^{f,b_1b_1}(s_2) = 10 \frac{9f_1f_2 - 216f_1 - 429f_2 + 346}{-9f_1 + 9f_2 + 40}$$

Figure 7: Zones \mathcal{Z}_{ij} for the SSSE and their boundaries in Example 1

$$\begin{aligned}
V_B^{f,b_1b_2}(s_1) &= 20 \frac{180 f_1 f_2 - 235 f_1 + 144 f_2 - 55}{-9 f_1 - 36 f_2 + 67} & V_B^{f,b_1b_2}(s_2) &= 20 \frac{180 f_1 f_2 - 207 f_1 + 176 f_2 - 83}{-9 f_1 - 36 f_2 + 67} \\
V_B^{f,b_2b_1}(s_1) &= 20 \frac{225 f_1 f_2 - 245 f_1 + 18 f_2 + 26}{-36 f_1 - 9 f_2 - 13} & V_B^{f,b_2b_1}(s_2) &= 20 \frac{225 f_1 f_2 - 225 f_1 + 48 f_2 - 2}{-36 f_1 - 9 f_2 - 13} \\
V_B^{f,b_2b_2}(s_1) &= 20 \frac{27 f_1 f_2 + 5 f_1 + 18 f_2 - 20}{9 f_1 - 9 f_2 + 10} & V_B^{f,b_2b_2}(s_2) &= 20 \frac{27 f_1 f_2 + 26 f_2 - 23}{9 f_1 - 9 f_2 + 10}.
\end{aligned}$$

C.1.2.2 Analysis of values. The objective is to compute, for every state, the best strategy $f \in W_A$ for Player A: the set $R_A(s)$. A SSSE will exist if and only if the intersection of these sets is nonempty.

C.1.2.2.1 Optimization for Player B. First of all, we identify the sets $R_B(f)$, that is, the solutions to Player B's MDP problem. This is done by calculating the “zones” where some policy $g = (b_i, b_j) \in W_B$ is optimal, formally defined as $\mathcal{Z}_{ij} = \{f \in W_A, (b_i, b_j) \in R_B(f)\}$.

We successively compute the differences $V_B^{f,b_1b_1}(s) - V_B^{f,b_1b_2}(s)$, $V_B^{f,b_1b_2}(s) - V_B^{f,b_2b_2}(s)$, etc. and we identify four critical lines separating the four zones \mathcal{Z}_{ij} $i, j = 1, 2$, as represented in Figure 7.

C.1.2.2.2 Optimization for Player A. According to the preferences of Player B, we have identified in the previous section a covering of W_A in four sets \mathcal{Z}_{ij} , on which we proceed to find Player A's maximum (this is a covering and not a partition, since the sets \mathcal{Z}_{ij} are not disjoint).

The plot of the functions $f \mapsto V_A^{f,\gamma_B(f)}(s)$ when $f \in W_A$, are displayed in Figure 8. We deduce from this that the global maximum of both $V_A^{f,\gamma_B(f)}(s_1)$ and $V_A^{f,\gamma_B(f)}(s_2)$ is

located at the point where all four zones meet. The zone that realizes the maximum is Z_{21} . The coordinates of this point are obtained e.g. by solving for (f_1, f_2) in the equations $V_B^{f, b_1 b_1}(s_1) = V_B^{f, b_1 b_2}(s_1) = V_B^{f, b_2 b_1}(s_1)$. The solution provides the value of f_2^* as the root of polynomial $p(f_2) = 3465 f_2^3 - 22604 f_2^2 + 26345 f_2 - 8516$ that belongs to $[0, 1]$, and f_1^* as a rational function of it. Finally, there exists a SSSE, given by the elements in Table 2.

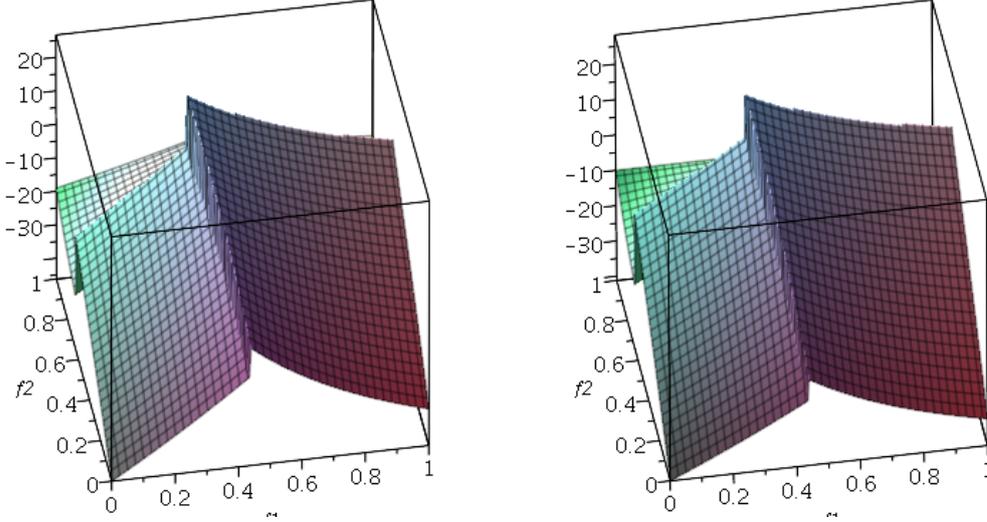


Figure 8: Value $V_A^{f, \gamma_B(f)}(s)$ in Example 1, for $s = s_1$ (left) and $s = s_2$ (right)

C.1.3 One-step Value Iteration from the SSSE

Since there is a natural candidate for a FPE which is the SSSE computed in Table 2, we can check whether this strategy satisfies the conditions for FPE. Replacing $v_B(s_1)$ and $v_B(s_2)$ by the numerical values of Table 2, we obtain the functions (up to rounding of floating-point values):

$$\begin{aligned} g_B(s_1, f, b_1) &= -14.25f + 3.13 & g_B(s_1, f, b_2) &= 11.01f - 5.62 \\ g_B(s_2, f, b_1) &= -15.25f + 9.13 & g_B(s_2, f, b_2) &= 17.01f - 11.62. \end{aligned}$$

Likewise for Player A, with the scrap value function $v_A(s_1) = 26.841$, $v_A(s_2) = 28.437$ (see again Table 2), the gains are given by:

$$\begin{aligned} g_A(s_1, f, b_1) &= 17.64f + 17.23 & g_A(s_1, f, b_2) &= -9.56f + 30.15 \\ g_A(s_2, f, b_1) &= 9.64f + 22.23 & g_A(s_2, f, b_2) &= -1.56f + 26.15. \end{aligned}$$

Finally, the decision problem of Player A is represented in Figure 9: for each state s , we plot the function $g_A(s, f)$:

$$f \mapsto g_A(s, f, b_1) \mathbf{1}_{\{g_B(s, f, b_1) \geq g_B(s, f, b_2)\}} + g_A(s, f, b_2) \mathbf{1}_{\{g_B(s, f, b_1) < g_B(s, f, b_2)\}}.$$

It is checked that the maximum is attained at the values of f that are given in Table 2. This means that the SSSE is a FPE.

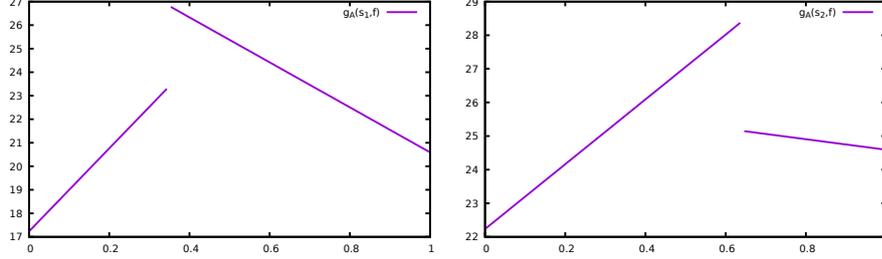


Figure 9: Value $g_A(s, f)$ in Example 1, for $s = s_1$ (left) and $s = s_2$ (right)

C.2 Analysis of Example 2

In this section, we provide details on the analysis of Example 2 introduced in Section 4.3. We compute the Stationary Strong Stackelberg Equilibria (SSSE) and discuss their existence. We compute the one-step operator and discuss its fixed points (FPE).

C.2.1 Data

The data is presented in Table 15 (also Table 3).

	b_1		b_2			b_1		b_2	
a_1	(1, 0)	(-1, -2)	(1, 0)	(-2, 1)	a_1	(1, 0)	(-1, -2)	(1, 0)	(-2, -2)
a_2	(1, 0)	(0, 0)	(1, 0)	(2, 0)	a_2	(0, 1)	(0, 1)	(1, 0)	(1, 1)
	State s_1					State s_2			

Table 15: Transition matrix and payoffs for each player in Example 2

Observe that state s_1 is absorbing whatever the players play, and state s_2 is followed by state s_1 , unless players play the combination (a_2, b_1) . This is represented in Figure 2. Also, in state s_2 , Player B's actions are indifferent to her.

C.2.2 Computation of the SSSE

C.2.2.1 Values of stationary strategies. Given a stationary strategy (f, g) , let us compute the value V_i^{fg} of this strategy for each player $i = A, B$. The notation is $f_i = f(s_i, a_1)$ and $g_i = g(s_i, b_1)$. We have:

$$\begin{aligned}
 v_A(s_1) &= f_1 g_1 [-1 + \beta v_A(s_1)] + f_1 (1 - g_1) [-2 + \beta v_A(s_1)] \\
 &\quad + (1 - f_1) g_1 [0 + \beta v_A(s_1)] + (1 - f_1) (1 - g_1) [2 + \beta v_A(s_1)]
 \end{aligned}$$

$$= \frac{1}{1 - \beta_A} (3f_1g_1 + 2 - 4f_1 - 2g_1) \quad (\text{C.1})$$

$$\begin{aligned} v_B(s_1) &= f_1g_1[-2 + \beta_B v_B(s_1)] + f_1(1 - g_1)[1 + \beta_B v_B(s_1)] \\ &\quad + (1 - f_1)g_1[0 + \beta_B v_B(s_1)] + (1 - f_1)(1 - g_1)[0 + \beta_B v_B(s_1)] \\ &= \frac{f_1(1 - 3g_1)}{1 - \beta_B}. \end{aligned} \quad (\text{C.2})$$

$$\begin{aligned} v_A(s_2) &= f_2g_2[-1 + \beta_A v_A(s_1)] + f_2(1 - g_2)[-2 + \beta_A v_A(s_1)] \\ &\quad + (1 - f_2)g_2[0 + \beta_A v_A(s_2)] + (1 - f_2)(1 - g_2)[1 + \beta_A v_A(s_1)] \\ &= \frac{1}{1 - (1 - f_2)g_2\beta_A} (2f_2g_2 - 3f_2 - g_2 + 1 + (1 - (1 - f_2)g_2)\beta_A v_A(s_1)) \end{aligned} \quad (\text{C.3})$$

$$\begin{aligned} v_B(s_2) &= f_2g_2[-2 + \beta_B v_B(s_1)] + f_2(1 - g_2)[-2 + \beta_B v_B(s_1)] \\ &\quad + (1 - f_2)g_2[1 + \beta_B v_B(s_2)] + (1 - f_2)(1 - g_2)[1 + \beta_B v_B(s_1)] \\ &= \frac{1}{1 - (1 - f_2)g_2\beta_B} (1 - 3f_2 + (1 - (1 - f_2)g_2)\beta_B v_B(s_1)). \end{aligned} \quad (\text{C.4})$$

C.2.2.2 Analysis of values. The objective is to compute, for every state, the best strategy $f \in W_A$ for Player A: the set $R_A(s)$. A SSSE will exist if and only if the intersection of these sets is nonempty. We start with the absorbing state s_1 and proceed with s_2 .

In addition to the existing notation, introduce:

$$R_B(s, f) = \{g \in W_B \mid V_B^{fg}(s) \geq V_B^{fh}(s), \forall h \in W_B\} \cap \prod_{s \in \{s_1, s_2\}} \{0, 1\}^{\mathcal{B}_s},$$

for the set of Player B's best reactions to some fixed strategy $f \in W_A$, when in state s . The computation will have the following steps. First, we compute the sets $R_B(s, f)$, $s \in \{s_1, s_2\}$ for individual states. We deduce $R_B(f)$ by intersecting them.

Next, we look for Player A's optimum: the maximum of $V_A^{fg}(s)$ where g is Player B's (global) response to strategy f . For state s_1 , which is absorbing, this turns out to be straightforward. For state s_2 , we focus the analysis on the sets $\mathcal{Z}_{ij} = \{f \in W_A \mid (b_i, b_j) \in R_B(s_2, f)\}$, that is, the subsets of W_A where playing (b_i, b_j) is optimal for Player B in state s_2 . We first identify these sets. In a second step, the maximum of Player A's gain is computed on each \mathcal{Z}_{ij} . Finally, the maximum of these maxima is identified.

Observe that the sets \mathcal{Z}_{ij} are not disjoint in general: when they intersect, Player B is indifferent between at least two (pure) strategies. She will then select a strategy that maximizes Player A's gain: the set $SR_B(f)$. However, the analysis does not need to identify precisely this best reaction of Player B.

C.2.2.2.1 Player B's response for state s_1 . We first look at the maximum of $V_B^{fg}(s_1) = v_B(s_1)$ with respect to (g_1, g_2) . Because state s_1 is absorbing, the expression of $v_i(s_1)$, $i = A, B$, depends only on f_1 and g_1 . The set $R_B(f, s_1)$ of Player

B's strategies maximizing $v_B(s_1)$, is of the form $R_B(f, s_1) = R_1 \times \{b_1, b_2\}$, where $R_1 = \arg \max\{v_B(s_1), g \in \{b_1, b_2\}\}$.

Clearly with (C.2), Player B prefers b_2 ($g_1 = 0$) when $f_1 \neq 0$. When $f_1 = 0$, she is indifferent between b_1 and b_2 . In other words: if $f_1 = 0$, then $R_1 = \{b_1, b_2\}$; if $f_1 \neq 0$, then $R_1 = \{b_2\}$.

C.2.2.2 Player B's response for state s_2 . We look now at the maximum of $V_B^{fg}(s_2)$ with respect to (g_1, g_2) , the values of (f_1, f_2) being fixed.

Preliminary. First we replace in (C.4) the term $v_B(s_1)$ by its value from (C.2):

$$\begin{aligned} v_B(s_2) &= \frac{1}{1 - (1 - f_2)g_2\beta_B} \left(1 - 3f_2 + (1 - (1 - f_2)g_2)\beta_B \frac{f_1(1 - 3g_1)}{1 - \beta_B} \right) \\ &= \frac{1 - f_1 - 3f_2 + 3f_1g_1}{1 - (1 - f_2)g_2\beta_B} + f_1 \frac{1 - 3g_1}{1 - \beta_B}. \end{aligned} \quad (\text{C.5})$$

We deduce:

$$\frac{\partial}{\partial g_1} v_B(s_2) = - \frac{3\beta_B f_1 (1 - (1 - f_2)g_2)}{(1 - \beta_B)(1 - (1 - f_2)g_2\beta_B)}.$$

As a function of (g_1, g_2) , $v_B(s_2)$ is strictly decreasing with respect to g_1 , when $\beta_B \neq 0$, $f_1 \neq 0$ and $1 - (1 - f_2)g_2 \neq 0$. Its maximum is then attained at $g_1 = 0$ for all values of g_2 . It is independent on g_1 in the other cases. We will analyze successively the cases $f_1 = 0$, and $f_1 > 0$.

We first assume $\beta_B \neq 0$. The case $\beta_B = 0$ will be handled separately.

The case $f_1 = 0$. Letting $f_1 = 0$ in the value of Player B, we arrive at:

$$v_B(s_2) = \frac{1 - 3f_2}{1 - (1 - f_2)g_2\beta_B},$$

which is independent of g_1 as determined above. Since $\beta_B > 0$, as a function of g_2 , this is strictly increasing if $1 - 3f_2 > 0$ and $f_2 \neq 1$, strictly decreasing if $1 - 3f_2 < 0$ and $f_2 \neq 1$, and constant if $1 - 3f_2 = 0$ or $f_2 = 1$. The optimal choice of g_2 is respectively 1 (B plays b_1), 0 (B plays b_2) and indifferent. In other words: if $f_2 < 1/3$, the strategy (f_1, f_2) belongs to $\mathcal{Z}_{11} \cap \mathcal{Z}_{21}$, if $1 > f_2 > 1/3$, it belongs to $\mathcal{Z}_{12} \cap \mathcal{Z}_{22}$, and if $f_2 = 1/3$ or $f_2 = 1$, it belongs to $\mathcal{Z}_{11} \cap \mathcal{Z}_{21} \cap \mathcal{Z}_{12} \cap \mathcal{Z}_{22}$. Player B's value is the same whatever she plays in these last cases.

The case $f_1 > 0$ and $f_2 > 0$. When $f_2 > 0$, $1 - (1 - f_2)g_2 \neq 0$. As we concluded in the preliminaries, $g_1 = 0$ in this case (Player B plays b_2 in state s_1). Points (f_1, f_2) with $f_2 > 0$ belong to either \mathcal{Z}_{21} or \mathcal{Z}_{22} , or both. Setting $g_1 = 0$ in (C.5), we have:

$$v_B(s_2) = \frac{1 - f_1 - 3f_2}{1 - (1 - f_2)g_2\beta_B} + \frac{f_1}{1 - \beta_B}.$$

The derivative of this function of g_2 is:

$$\frac{\partial}{\partial g_2} v_B(s_2) = \frac{(1 - f_1 - 3f_2)\beta_B(1 - f_2)}{(1 - (1 - f_2)g_2\beta_B)^2} \quad (\text{C.6})$$

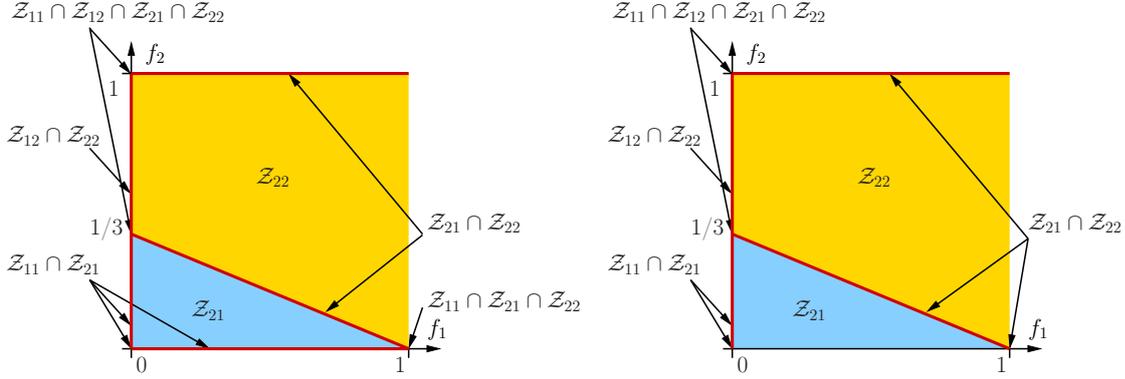


Figure 10: Zones of W_A for the response sets of Player B, $R_B(f, s_2)$ (left) and $R_B(f)$ (right); red lines are points that belong to several zones

and is either 0 if $f_2 = 1$, or else has the sign of $1 - f_1 - 3f_2$. In the cases $f_2 = 1$ and $f_1 = 1 - 3f_2$, Player B is therefore indifferent and the strategy f belongs to $Z_{21} \cap Z_{22}$. If $f_1 < 1 - 3f_2$, Player B will prefer $g_2 = 1$ (play b_1): those points are in Z_{21} . If $f_1 > 1 - 3f_2$, Player B will prefer $g_2 = 0$ (play b_2): those points are in Z_{22} . The situation is depicted in Figure 10.

The case $f_1 > 0$ and $f_2 = 0$. Setting $f_2 = 0$ in (C.5), we arrive at:

$$\begin{aligned} v_B(s_2) &= \frac{1 - f_1 + 3f_1g_1}{1 - \beta_B g_2} + f_1 \frac{1 - 3g_1}{1 - \beta_B} \\ &= \frac{1}{1 - \beta_B} - \frac{\beta_B}{1 - \beta_B} \frac{(1 - g_2)(3f_1g_1 - f_1 + 1)}{1 - \beta_B g_2}. \end{aligned}$$

It is seen that the second term in the right-hand side is always positive, and is zero when either $g_2 = 1$, or $f_1 = 1$ and $g_1 = 0$. The maximum possible for Player B's value is therefore $v_B(s_2) = 1/(1 - \beta_B)$, and it is attained when either $g_2 = 1$ (independently of f_1 and g_1 : in that case, the point is in $Z_{11} \cap Z_{21}$), or when $f_1 = 1$ and $g_1 = 0$ (independently of g_2 : in that case, the point is in Z_{22} as well).

C.2.2.2.3 Optimization for Player A. Preliminary. The best response of Player B is obtained as $R_B(s_1, f) \cap R_B(s_2, f)$ where these sets have been computed in Section C.2.2.2.2. When $f_1 > 0$, $R_B(s_1, f) = \{b_2\} \times \{b_1, b_2\}$ so that the constraint imposed on $R_B(s_2, f)$ is that $g_1 = b_2$. When $f_1 = 0$, no constraint is imposed. The set $R_B(f)$ is then as depicted in Figure 10 (right). The difference with $R_B(f, s_2)$ is that the line $\{(f_1, 0), 0 < f_1 \leq 1\}$ does not belong to Z_{11} .

The value $v_A(s_1)$ is given in (C.1). When we replace the term $v_A(s_1)$ by this value in (C.3), we get:

$$v_A(s_2) = \frac{1}{1 - (1 - f_2)g_2\beta_A} \left(2f_2g_2 - 3f_2 - g_2 + 1 \right)$$

$$+ (1 - (1 - f_2)g_2) \frac{\beta_A}{1 - \beta_A} (3f_1g_1 + 2 - 4f_1 - 2g_1) \Big). \quad (\text{C.7})$$

State s_1 . Similarly as for Player B's optimization, Player A's problem does not depend on f_2 and reduces to the best choice of f_1 . Assume first that $f_1 > 0$. The set $SR_B(f)$ is not known yet, but whatever it is, Player B plays b_2 ($g_1 = 0$). Player A's value is then $(2 - 4f_1)/(1 - \beta_A)$ which would be maximized with $f_1 = 0$ (but the value is outside the case studied).

If $f_1 = 0$, Player A, has gain, from (C.1): $(2 - 2g_1)/(1 - \beta_A)$. In that case also, B will play $g_1 = 0$.

So Player A gets to choose f_1 knowing that $g_1 = 0$ in (C.1). The optimum is reached at $f_1 = 0$. The best for Player A is when she plays a_2 and B reacts with b_2 . In other terms: $R_A(s_1) = \{a_2\} \times \mathbb{P}(A_{s_2})$. The values are then:

$$v_A^*(s_1) = \frac{2}{1 - \beta_A} \quad v_B^*(s_1) = 0. \quad (\text{C.8})$$

State s_2 . According to the preferences of Player B, we have identified the covering of W_A by zones \mathcal{Z}_{ij} , $i, j = 1, 2$, on which we proceed to find Player A's maximum. See again Figure 10 for a visualization. Let $m_{ij} = \max\{v_A(s_2) | (f_1, f_2) \in \mathcal{Z}_{ij}\}$.

Zone \mathcal{Z}_{11} : \mathcal{Z}_{11} is the segment $\{(0, f_2), f_2 \in [0, 1/3]\}$.¹ Since $g_1 = g_2 = 1$, we have from (C.7):

$$v_A(s_2) = -\frac{f_2(1 - \beta_A(1 - f_1))}{1 - (1 - f_2)\beta_A} \frac{1}{1 - \beta_A}.$$

All the factors being positive, it is seen that $v_A(s_2) \leq 0$. If $f_2 = 0$, $v_A(s_2) = 0$ and the maximum is realized. As soon as $f_2 > 0$, $v_A(s_2) < 0$. In conclusion, $m_{11} = 0$, a value realized at point $(0, 0)$.

Zone \mathcal{Z}_{12} : \mathcal{Z}_{12} is the segment $\{(0, f_2), f_2 \in [1/3, 1]\}$. Player A's value is:

$$v_A(s_2) = 1 - 3f_2.$$

Its maximum is therefore attained at $f_2 = 1/3$, where $m_{12} = 0$.

Zone \mathcal{Z}_{21} : \mathcal{Z}_{21} is the region $\{(f_1, f_2) \in W_A | 0 \leq f_2 \leq (1 - f_1)/3\} \cup \{(f_1, 1) | 0 \leq f_1 \leq 1\}$. Player A's value writes as:

$$v_A(s_2) = \frac{f_2}{1 - (1 - f_2)\beta_A} \frac{\beta_A(3 - 4f_1) - 1}{1 - \beta_A}.$$

If $f_2 = 0$, this is constant equal to 0. If $f_2 > 0$, this is decreasing with respect to f_1 (strictly if $\beta_A > 0$) and therefore maximal at $f_1 = 0$. The value obtained reduces to

$$\frac{f_2}{1 - \beta_A(1 - f_2)} \frac{3\beta_A - 1}{1 - \beta_A}$$

¹As observed in Section C.2.2.2.3, for $f = (f_1, 0)$, $f_1 > 0$, strategy (b_1, b_1) belongs to $R_B(s_2, f)$ but not to $R_B(f)$.

and its behavior as a function of f_2 depends on the sign of $3\beta_A - 1$. If $\beta_A \leq 1/3$, the maximum is at $f_2 = 0$ with value 0. If $\beta_A > 1/3$, the maximum is at $f_2 = 1$. We have the following table for locating the maximum m_{21} :

case	m_{21}	location
$0 \leq \beta_A < 1/3$	0	$f_1 \in [0, 1], f_2 = 0$
$\beta_A = 1/3$	0	$f_1 = 0, f_2 \in [0, 1/3]$
$\beta_A > 1/3$	$\frac{3\beta_A - 1}{1 - \beta_A}$	$f_1 = 0, f_2 = 1$

Zone \mathcal{Z}_{22} : \mathcal{Z}_{22} is the region $\{(f_1, f_2) \in W_A | (1 - f_1)/3 \leq f_2 \leq 1\}$. Player A's value is:

$$v_A(s_2) = 1 - 3f_2 + \frac{2\beta_A(1 - 2f_1)}{1 - \beta_A}.$$

It is therefore maximized when f_2 is as small as possible, which occurs when $f_1 = 1 - 3f_2$ for all values of f_1 . When $f_1 = 1 - 3f_2$, the value of Player A is given by:

$$3f_2 \frac{5\beta_A - 1}{1 - \beta_A} + \frac{1 - 3\beta_A}{1 - \beta_A} = f_1 \frac{1 - 5\beta_A}{1 - \beta_A} + \frac{2\beta_A}{1 - \beta_A}.$$

The location of m_{22} depends on the sign of $5\beta_A - 1$. Enumerating the cases, we have the following table:

case	m_{22}	location
$0 \leq \beta_A < 1/5$	$\frac{1 - 3\beta_A}{1 - \beta_A}$	$f_1 = 1, f_2 = 0$
$\beta_A = 1/5$	$\frac{1}{2}$	$f_1 = 1 - 3f_2, f_2 \in [0, 1/3]$
$\beta_A > 1/5$	$\frac{2\beta_A}{1 - \beta_A}$	$f_1 = 0, f_2 = 1/3$

The case $\beta_B = 0$. If $\beta_B = 0$, Player B's value in state s_2 is $1 - 3f_2$ and does not depend on g_1 nor on g_2 . In other terms, $\mathcal{Z}_{11} = \mathcal{Z}_{12} = \mathcal{Z}_{21} = \mathcal{Z}_{22} = W_A$. Player A's maximization problem is not constrained, and consists in choosing all four variables f_1, f_2, g_1, g_2 in (C.7). As in the general analysis, we have a dichotomy: either $\beta_A(1 - (1 - f_2)g_2) = 0$, either it is not 0.

In the first case, the maximization does not involve f_1 or g_1 . Two possibilities:

- either $\beta_A = 0$: Player A's value is: $v_A(s_2) = 2f_2g_2 - 3f_2 - g_2 + 1$ which is maximized at $f_2 = g_2 = 0$. The gain of Player A is then 1.
- or $f_2 = 0$ and $g_2 = 1$: in that situation, Player A's value is always 0.

In the second case, one starts with maximizing with respect to (f_1, g_1) , which yields $f_1 = g_1 = 0$ because it is the same as playing the static game in s_1 . Then the analysis proceeds to show that $f_2 = g_2 = 0$ realizes the maximum.

The conclusion is that in all cases of β_A , the optimal strategy of Player A is to play $f_1 = f_2 = 0$ (play a_2 in both states). Player B will react with $g_2 = 0$ and $g_1 = 0$ (unless $\beta_A = 0$, in which case the value of g_1 is irrelevant for Player A). The values obtained are:

$$v_A^*(s_2) = \frac{1 + \beta_A}{1 - \beta_A} \quad v_B^*(s_2) = 1. \quad (\text{C.9})$$

C.2.2.3 Synthesis: existence and value of SSSE Summing up the results obtained in Section C.2.2.2.3 we have four cases:

$0 \leq \beta_A < 1/5$ **and** $\beta_B > 0$: The sets $R_A(s)$ are given by the following tables:

$R_A(s_1)$	s_1	s_2	$R_A(s_2)$	s_1	s_2
Play of A	a_2	$\mathbb{P}(\mathcal{A}_{s_2})$	Play of A	a_1	a_2
Play of B	b_2	$\{b_1, b_2\}$	Play of B	b_2	b_2
$v_A(s_1)$	$\frac{2}{1 - \beta_A}$		$v_A(s_2)$	$\frac{1 - 3\beta_A}{1 - \beta_A}$	
$v_B(s_1)$	0		$v_B(s_2)$	$\frac{1}{1 - \beta_B}$	

The intersection $R_A(s_1) \cap R_A(s_2)$ is empty: there exist no SSSE.

$\beta_A = 1/5$ **and** $\beta_B > 0$: The set $R_A(s_1)$ is as above. The set $R_A(s_2)$ is described in the following table:

$R_A(s_2)$	s_1	s_2
Play of A	$(1 - 3\phi, 3\phi)$	$(\phi, 1 - \phi)$
Play of B	b_2	$\{b_1, b_2\}$
$v_A(s_2)$	1/2	
$v_B(s_2)$	$\frac{1 - 3\phi}{1 - \beta_B}$	

where $\phi \in [0, 1/3]$. When $\phi = 1/3$, the corresponding element of $R_A(s_2)$ is equal to an element of $R_A(s_1)$. There is therefore an intersection, which is a SSSE, namely: $f_1 = 0, f_2 = 1/3$.

$\beta_A > 1/5$ **and** $\beta_B > 0$: The set $R_A(s_1)$ is still as above and the set $R_A(s_2)$ is the element described in the following table:

$R_A(s_2)$	s_1	s_2
Play of A	a_2	$(1/3, 2/3)$
Play of B	b_2	b_2
$v_A(s_2)$	$\frac{2\beta_A}{1 - \beta_A}$	
$v_B(s_2)$	0	

This set is included in $R_A(s_1)$. There exists then a unique SSSE equal to this element.

$\beta_B = 0$: The set $R_A(s_1)$ is still as above and the set $R_A(s_2)$ is the element described in the following table:

$R_A(s_2)$	s_1	s_2
Play of A	a_2	a_2
Play of B	b_2	b_2
$v_A(s_2)$	$\frac{1 + \beta_A}{1 - \beta_A}$	
$v_B(s_2)$	1	

This set is included in $R_A(s_1)$. There exists then a unique SSSE equal to this element.

C.2.2.4 Concluding comments on SSSE. The non-existence of SSSE for $\beta_A < \frac{1}{5}$ and $\beta_B > 0$ comes from the fact that $R_A(s_1) \cap R_A(s_2) = \emptyset$ in this case. The leader plays different strategies if she starts in state s_1 than if she starts in s_2 . We explain briefly the case for $\beta_A = 0$.

Note that if the game starts in s_1 , then the optimal policy for the leader, will be the same than in the static case, that is a game with only state s_1 . The set of best strategies for the leader is to play $f_1 = 0$ (play a_2) in state s_1 and any arbitrary $f_2 \in [0, 1]$. The follower reacts with b_2 in state s_1 and an arbitrary response in s_2 . In this case, the value functions are $v_A(s_1) = 2$ and $v_B(s_1) = 0$.

If the game starts in s_2 , the value of Player A is, as we have already seen: $v_A(s_2) = 2f_2g_2 - 3f_2 - g_2 + 1$. It is maximized when $f_2 = g_2 = 0$, that is, when (a_2, b_2) is played. If Player A performs a_2 , Player B is indifferent with immediate rewards (she gains 1 whatever she plays) and she has to arbitrate between staying forever in state s_2 , and receive an additional $\beta_B/(1 - \beta_B)$ by playing b_1 , or jumping to state s_1 and gain $\beta_B v_B(s_1)$ by playing b_2 . Player A prefers this second situation because of the immediate reward, not because of future gains. Her interest is therefore to use the tie-breaking rule of Player B to induce her to play b_2 and not b_1 in state s_2 .

From (C.2) and (C.4), and assuming $f_2 = 0$, Player B's value is:

$$v_B(s_2) = \frac{1 - f_1 + 3f_1g_1}{1 - g_2\beta_B} + f_1 \frac{1 - 3g_1}{1 - \beta_B}$$

$$= \frac{1}{1 - \beta_B} + \beta_B(g_2 - 1) \frac{3f_1g_1 - f_1 + 1}{(1 - \beta_B)(1 - g_2\beta_B)}.$$

As it turns out, $3f_1g_1 - f_1 + 1 \geq 0$ and it can be $= 0$ only when $f_1 = 1$ and $g_1 = 0$. If Player A plays any $f_1 \neq 0$, the only way Player B can maximize her value is to play $g_2 = 1$ which is contrary to Player A's preference. On the other hand, if Player A plays $f_1 = 1$, Player B has can choose either $g_1 = 0$, or $g_2 = 1$ (or both). This tie is broken in favor of A, whose value is, when $f_1 = 1$ and $f_2 = 0$: $v_A(s_2) = 1 - g_2$. Therefore B picks $g_2 = 0$ (g_1 being indifferent), thereby fulfilling Player A's objective. Playing $f_1 = 1$ (action a_2) in state s_1 is however not optimal for Player A since action a_2 strictly dominates it. This explains the lack of SSSE.

This argument does not work when $\beta_B = 0$. In that case, Player B is indifferent and spontaneously breaks the tie in favor of Player A by playing b_2 in state s_2 ,

Observe here in passing that the ‘‘myopic’’ case does not always bring a simplification to the problem. Despite the fact that the myopic follower condition guarantees the existence of FPE and SSSE, the myopic leader condition does not guarantee the existence of SSSE.

C.2.3 The one-step Value Iteration operator

C.2.3.1 Computation of the operator. We compute here the operator T , which maps pairs of value functions to pairs of value functions. We restrict our attention to those pairs where the value at state s_1 is the fixed-point obtained in (C.8). There are therefore two remaining variables: $v_A(s_2)$ and $v_B(s_2)$. In order to clarify formulas a bit, we use the symbols $w := v_A(s_2)$ and $z := v_B(s_2)$. The images $(Tv)_A(s_2)$ and $(Tv)_B(s_2)$ will be denoted with w' and z' respectively.

By using the values in (C.8), we arrive at a single-state game parametrized by the values of w and z :

	b_1	b_2
a_1	$\frac{3\beta_A - 1}{1 - \beta_A}, -2$	$\frac{4\beta_A - 2}{1 - \beta_A}, -2$
a_2	$\beta_A w, 1 + \beta_B z$	$\frac{1 + \beta_A}{1 - \beta_A}, 1$

C.2.3.1.1 Operators T_i^{fg} First of all, we write down the expressions for $(T_i^{fg}v)(s_2)$, $i = 1, 2$ for arbitrary strategies $f \in W_A$, $g \in W_B$. Those are deduced from the computation in Section C.2.2.1, by adapting the right-hand side of equations, and using values in (C.8). Then,

$$(T^{fg}v)_A(s_2) = 2f_2g_2 - 3f_2 - g_2 + 1 + (1 - f_2)g_2\beta_A w + (1 - (1 - f_2)g_2) \frac{2\beta_A}{1 - \beta_A} \quad (\text{C.10})$$

$$(T^{fgv})_B(s_2) = 1 - 3f_2 + (1 - f_2)g_2\beta_B z . \quad (\text{C.11})$$

C.2.3.1.2 Optimization of Player B. We now identify the reaction sets of Player B to a given strategy of Player A, a given scrap value and state s_2 : the sets $R_B(s_2, f, v_B)$.

From (C.11) we see that Player B's optimal reaction to a given f_2 depends on the sign of z , unless $f_2 = 1$ or $\beta_B = 0$. We assume $\beta_B > 0$ in the remainder of this section. The case $\beta_B = 0$ will be studied in Section C.2.3.3.

If $f_2 = 1$, or $f_2 < 1$ and $z = 0$, then Player B is indifferent: $R_B(s_2, f, v_B) = \{b_1, b_2\}$.

If $f_2 < 1$ and $z > 0$, $R_B(s_2, f, v_B) = \{b_1\}$. If $f_2 < 1$ and $z < 0$, $R_B(s_2, f, v_B) = \{b_2\}$.

C.2.3.1.3 Optimization of Player A. Player A's reward when Player B plays b_1 ($g_2 = 1$) is:

$$v_A(s_2) = -f_2 + f_2 \frac{2\beta_A}{1 - \beta_A} + (1 - f_2)\beta_A w = f_2 \frac{3\beta_A - 1}{1 - \beta_A} + (1 - f_2)\beta_A w.$$

This is the situation when $z > 0$. The maximum of this expression with respect to f_2 depends on the location of w with respect to $(3\beta_A - 1)/\beta_A/(1 - \beta_A)$. Player A's reward when Player B plays b_2 ($g_2 = 0$) is:

$$v_A(s_2) = 1 - 3f_2 + \frac{2\beta_A}{1 - \beta_A} .$$

This is the situation when $z < 0$. The maximum of this expression with respect to f_2 is when $f_2 = 0$. When $z = 0$, these two values must be compared. Which one is the largest depends on the location of w with respect to $(1 + \beta_A)/\beta_A/(1 - \beta_A)$.

After optimization with respect to f_2 , we have the following cases:

Case $z < 0$: here, Player B always plays b_2 and A plays a_2 . We have the fixed mapping:

$$(w', z') = \left(\frac{1 + \beta_A}{1 - \beta_A}, 1 \right) =: P_{22} .$$

Case $z > 0$: here, Player B always plays b_1 .

Case $w < (3\beta_A - 1)/\beta_A/(1 - \beta_A)$: A plays a_1 . The mapping is also fixed:

$$(w', z') = \left(\frac{3\beta_A - 1}{1 - \beta_A}, -2 \right) =: P_{11} .$$

Case $w > (3\beta_A - 1)/\beta_A/(1 - \beta_A)$: A plays a_2 . The mapping is linear:

$$(w', z') = (\beta_A w, 1 + \beta_B z) =: T_{21}(w, z) .$$

Case $w = (3\beta_A - 1)/\beta_A/(1 - \beta_A)$: A plays anything she wants. Let ϕ be this strategy, the mapping is:

$$(w', z') = (\beta_A w, 1 - 3\phi + \beta_B(1 - \phi)z) =: T_{\phi,0}(w, z) . \quad (\text{C.12})$$

Case $z = 0$: Player B is indifferent with her own reward. The tie is broken in favor of Player A.

Case $w \leq (1 + \beta_A)/\beta_A/(1 - \beta_A)$: A prefers (a_2, b_2) . The mapping is to P_{22} .

Case $w > (1 + \beta_A)/\beta_A/(1 - \beta_A)$: A prefers (a_2, b_1) . The mapping is $T_{21}(w, z)$.

In a more compact form, the sub-operator for this game in state s_2 is as follows:

$$(w', z') = \begin{cases} P_{22} = \left(\frac{1+\beta_A}{1-\beta_A}, 1\right) & \text{if } z < 0 \\ P_{11} = \left(\frac{3\beta_A-1}{1-\beta_A}, -2\right) & \text{if } w < \frac{3\beta_A-1}{\beta_A(1-\beta_A)} \text{ and } z > 0 \\ T_{21}(w, z) = (\beta_A w, 1 + \beta_B z) & \text{if } w > \frac{3\beta_A-1}{\beta_A(1-\beta_A)} \text{ and } z > 0 \\ (\beta_A w, 1 - 3f_2 + (1 - f_2)\beta_B z) & \text{if } w = \frac{3\beta_A-1}{\beta_A(1-\beta_A)} \text{ and } z > 0 \\ T_{21}(w, 0) = (\beta_A w, 1) & \text{if } w > \frac{1+\beta_A}{\beta_A(1-\beta_A)} \text{ and } z = 0 \\ P_{22} = \left(\frac{1+\beta_A}{1-\beta_A}, 1\right) & \text{if } w \leq \frac{1+\beta_A}{\beta_A(1-\beta_A)} \text{ and } z = 0. \end{cases} \quad (\text{C.13})$$

In summary, the plane (w, z) is partitioned in three zones, as illustrated in Figure 11. The three zones are represented: \mathcal{Z}_{11} in green: points of this zone map to P_{11} ; \mathcal{Z}_{22} in blue: points of this zone map to P_{22} ; \mathcal{Z}_{21} in yellow: points of this zone are mapped by operator T_{21} . This mapping T_{21} admits as fixed point:

$$P_{21} := \left(0, \frac{1}{1 - \beta_B}\right).$$

In all cases, P_{22} is in the zone \mathcal{Z}_{21} and point P_{11} is in the zone \mathcal{Z}_{22} .

Other notable points are: Point C , the triple point and Point $L = ((1 + \beta_A)/\beta_A/(1 - \beta_A), 0)$ which delimits the two cases when $z = 0$: the horizontal segment $(-\infty, L]$ belongs to zone \mathcal{Z}_{22} and the segment (L, ∞) belongs to zone \mathcal{Z}_{21} .

Point S is the candidate SSSE with coordinates $(2\beta_A/(1 - \beta_A), 0)$. Points C and S are always located to the left of Point L , which means that they both belong to \mathcal{Z}_{22} .

The solid red line, that is, the half-line $\{w = (3\beta_A - 1)/\beta_A/(1 - \beta_A); z > 0\}$, is the frontier between \mathcal{Z}_{21} and \mathcal{Z}_{22} . It has a special status: by picking any number ϕ in (C.12), an infinite number of mappings $T_{\phi,0}$ can be chosen. Its image is represented as the dashed red line.

C.2.3.2 Dynamics of the operator. Two main situations occur: either $\beta_A > 1/3$, either $\beta_A \leq 1/3$. Those are represented in Figure 11, on the left and the right, respectively.

C.2.3.2.1 Case $\beta_A > 1/3$. In this situation, the fixed-point of operator T_{21} is located in zone \mathcal{Z}_{11} . The iterations of T starting in \mathcal{Z}_{22} (with $z < 0$ or in the segment $(-\infty, L]$) will be mapped first to P_{22} , then to $T_{21}P_{22}$, then $T_{21}^2P_{22}$ etc. Eventually, the sequence $T_{21}^k P_{22}$ will enter zone \mathcal{Z}_{11} . Then the value will be mapped to P_{11} and repeat the cycle. This is in particular the case of the SSSE which belongs to \mathcal{Z}_{22} , see Section C.2.4.

Iterations of T starting in zone \mathcal{Z}_{11} will map to P_{11} and continue as above.

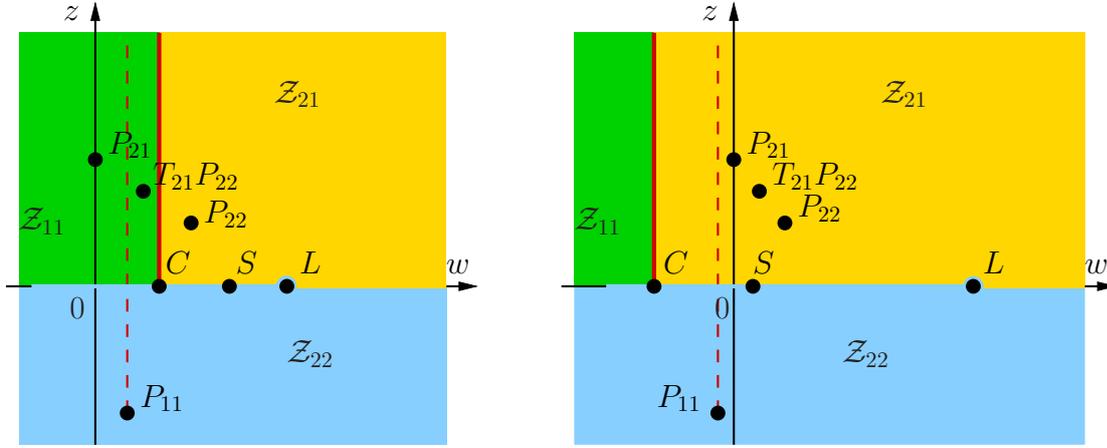


Figure 11: Zones of the mapping T : left, with $\beta_A > 1/3$; right, with $\beta_A < 1/3$

Iterations of T starting in zone Z_{21} will follow a sequence of points in the direction of P_{21} and eventually enter Z_{11} , then follow the cycle as above.

Depending on the choice of ϕ , the mapping from the frontier $Z_{11} - Z_{21}$ may lead to either Z_{11} or Z_{22} . It then follows the patterns described.

In all cases, there is a limit cycle. The length of this cycle can range from 3 to ∞ depending on the proximity of β_A to $1/3$.

C.2.3.2.2 Case $\beta_A \leq 1/3$. In this situation, the fixed-point P_{21} of operator T_{21} is located in zone Z_{21} : it is a FPE. The iterations of T starting in Z_{21} will converge to it.

Iterations of T starting in Z_{22} will be mapped first to P_{22} , then converge to P_{21} . Likewise, iterations starting in Z_{11} will be mapped first to P_{11} in Z_{22} , then converge to P_{21} . Iterations starting on the frontier $Z_{11} - Z_{21}$ either converge directly, or make one step in Z_{22} .

C.2.3.3 The operator when $\beta_B = 0$. In the special case $\beta_B = 0$, the expression for operator T_B^{fg} becomes, from *e.g.* (C.11):

$$(T^{fg}v)_B(s_2) = 1 - 3f_2 .$$

Since Player B is indifferent, Player A gets to optimize her own reward, namely,

$$(T^{fg}v)_A(s_2) = 2f_2g_2 - 3f_2 - g_2 + 1 + (1 - f_2)g_2\beta_A w + (1 - (1 - f_2)g_2) \frac{2\beta_A}{1 - \beta_A}$$

with respect to f_2 and g_2 . This optimization concludes that $f_2 = 0$ is always better. Finally, the expression for the operator is:

$$(Tv)_A(s_2) = \max \left\{ \beta_A w, \frac{1 + \beta_A}{1 - \beta_A} \right\}$$

$$(Tv)_B(s_2) = 1 .$$

This operator is contractive and has the function in (C.9) as fixed point.

C.2.4 One-step Value Iteration from the SSSE

In this section, we compute the action of VI on the SSSE computed in Section C.2.2.3, that is:

$$v_A^*(s_1) = \frac{2}{1 - \beta_A} \quad v_A^*(s_2) = \frac{2\beta_A}{1 - \beta_A} \quad v_B^*(s_1) = v_B^*(s_2) = 0. \quad (\text{C.14})$$

The purpose of this is to check again that the SSSE is not an FPE, but also to open the door to an investigation of nonstationary strategies in the context of Dynamic Stackelberg games.

Using the operator defined in Section C.2.3.1, we find that $z = 0$. We must compare $w = (2\beta_A)/(1 - \beta_A)$ to $L = (1 + \beta_A)/\beta_A/(1 - \beta_A)$ and find that $w \leq L$. Therefore the FPE lies in the zone \mathcal{Z}_{22} and it is mapped to P_{22} by T . The value obtained is then:

$$(Tv^*)_A(s_2) = \frac{1 + \beta_A}{1 - \beta_A} \quad (Tv^*)_B(s_2) = 1,$$

a function we have identified as the SSSE when $\beta_B = 0$, see (C.9).

Interpretation. Independently of β_A , the value of both players is improved by applying the operator once, starting from the SSSE. It means that the policy in which Player A announces she: 1) plays once a_2 ; 2) plays the SSSE afterwards, gives a better reward to both players.

This nonstationary policy actually *menaces* to play the SSSE, since at the first step players will play (a_2, b_2) and send the game to state s_1 . The mixed strategy $(1/3, 2/3)$ of the SSSE will never be played.

C.3 Analysis of Example 3

We provide here the technical justifications for the claims of Section 4.4 on Example 3.

C.3.1 Data

The data of this example is presented in Table 16 (see also Table 6):

		State s_1		State s_2	
		b_1	b_2	b_1	b_2
a_1		(1, 0)	(0, 1)	(0, 1)	(1, 0)
		(1, -1)	(0, 1)	(-1, 0)	(0, 1)
a_2		(0, 1)	(0, 1)	(1, 0)	(0, 1)
		(-1, 1)	(-1, -1)	(0, 1)	(1, -1)

Table 16: Transition matrix and payoffs for each player in Example 3

We claim that the following pair of strategies (f^*, g^*) and value functions (v_A^*, v_B^*) constitute *both* a SSSE and a FPE:

$$f^*(s_1, a_1) = 1 \quad f^*(s_2, a_1) = 5 - \sqrt{19} \quad v_A^*(s_1) = \frac{-3 + \sqrt{19}}{5} \quad v_A^*(s_2) = \frac{-6 + 2\sqrt{19}}{5} \quad (\text{C.15})$$

$$g^*(s_1) = b_2 \quad g^*(s_2) = b_2 \quad v_B^*(s_1) = \frac{16 - 2\sqrt{19}}{5} \quad v_B^*(s_2) = \frac{22 - 4\sqrt{19}}{5}. \quad (\text{C.16})$$

In order to support this claim, we need to construct the reaction functions $\gamma_B(f^*)$ (for the SSSE) and $\gamma_B(s, f^*, v^*)$ (for the FPE). As a common preliminary step, we first list the one-step rewards corresponding to a general strategy $f \in W_A$ and any possible action $g \in \mathcal{B}_s$, for all states s . The notation is $f_1 = f(s_1, a_2)$ and $f_2 = f(s_2, a_1)$.

$$\begin{aligned} h_A(s_1, f, b_1, v_A) &= 2f_1 - 1 + \beta_A[f_1 v_A(s_1) + (1 - f_1)v_A(s_2)] \\ h_A(s_1, f, b_2, v_A) &= f_1 - 1 + \beta_A v_A(s_2) \\ h_A(s_2, f, b_1, v_A) &= -f_2 + \beta_A[f_2 v_A(s_2) + (1 - f_2)v_A(s_1)] \\ h_A(s_2, f, b_2, v_A) &= 1 - f_2 + \beta_A[f_2 v_A(s_1) + (1 - f_2)v_A(s_2)] \\ h_B(s_1, f, b_1, v_B) &= 1 - 2f_1 + \beta_B[f_1 v_B(s_1) + (1 - f_1)v_B(s_2)] \\ h_B(s_1, f, b_2, v_B) &= -1 + 2f_1 + \beta_B v_B(s_2) \\ h_B(s_2, f, b_1, v_B) &= 1 - f_2 + \beta_B[f_2 v_B(s_2) + (1 - f_2)v_B(s_1)] \\ h_B(s_2, f, b_2, v_B) &= 2f_2 - 1 + \beta_B[f_2 v_B(s_1) + (1 - f_2)v_B(s_2)]. \end{aligned}$$

C.3.2 Computation of the SSSE

Given a strategy $f \in W_A$, we know from MDP theory that the best response of Player B is found among the four pure strategies of $W_B \times W_B$. The value of each of these strategies for both players is obtained by solving the four equations of the list above, where the combination of s_i and b_j is the one sought, and where $h_i(s, f, b, v_i)$ is replaced with $v_i(s_i)$.

When setting $\beta_B = 1/2$, the four values of $V_B^{fg}(s_2)$ turn out to be:

$$\begin{aligned} V_B^{f, b_1 b_1}(s_2) &= \frac{6(1 - f_2)(1 - f_1)}{3 - f_1 - f_2} & V_B^{f, b_1 b_2}(s_2) &= \frac{2(4f_1 f_2 - f_1 - 5f_2 + 2)}{f_1 - f_2 - 2} \\ V_B^{f, b_2 b_1}(s_2) &= \frac{2(1 - f_2)(2f_1 + 1)}{3 - f_2} & V_B^{f, b_2 b_2}(s_2) &= \frac{2(2f_1 f_2 + 3f_2 - 2)}{f_2 + 2}. \end{aligned}$$

The comparison of these four values determines four zones delimited by lines as in Figure 12.

When evaluating Player A's value on each of these zones, including the boundaries, it is found that A's optimum lies at the point specified in (C.15). This point lies at the boundary of B's best responses $g = (b_2, b_1)$ and (b_2, b_2) (line L6), where $f_1 = 1$. Player B breaks the tie in favor of Player A by playing b_2 .

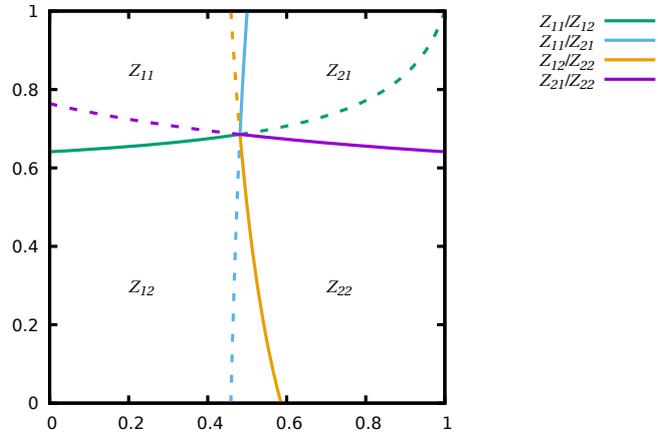


Figure 12: Best response zones for Example 3

C.3.3 Verification of the FPE

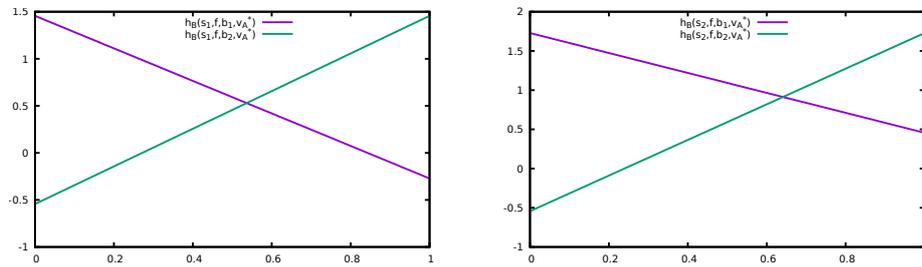
When replacing $v_B(s_1)$ and $v_B(s_2)$ by the values in (C.16), the functions $h_B(s, \cdot)$ of the variable $f(s, a_1)$ are as shown in Figure 13. The intersections occur for the values $f(s_1, a_1) = f_1^* := (23 + \sqrt{19})/51$ and $f(s_2, a_1) = f_2^* := f^*(s_2, a_1)$ as in (C.15). In addition, it is checked that for these values of $f(s, a_1)$,

$$h_A(s_1, f, b_1, v_A^*) > h(s_1, f, b_2, v_A^*) \quad h_A(s_2, f, b_1, v_A^*) < h(s_2, f, b_2, v_A^*)$$

so that the “strong” response of Player B is:

$$\gamma_B(s_1, f, v) = \begin{cases} b_1 & \text{if } f_1 \leq f_1^* \\ b_2 & \text{if } f_1 > f_1^* \end{cases} \quad \gamma_B(s_2, f, v) = \begin{cases} b_1 & \text{if } f_2 < f_2^* \\ b_2 & \text{if } f_2 \geq f_2^*. \end{cases}$$

In particular, this confirms that $\gamma_B(s, f^*, v) = \{b_2\}$ for $s = s_1, s_2$.

Figure 13: Player B's responses as a function of $f_1 = f(s, a_1)$ in Example 3; state s_1 (left) and state s_2 (right)

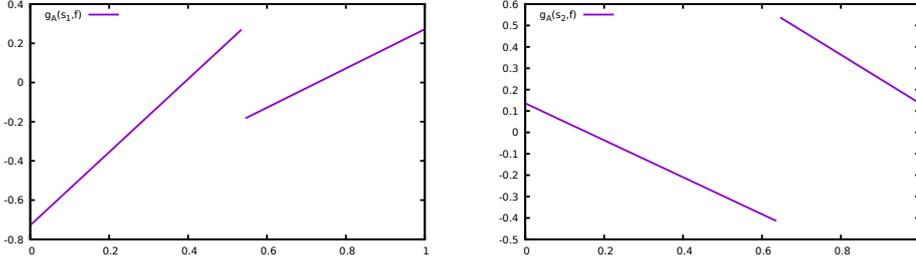


Figure 14: Value of Player A as a function of $f_1 = f(s, a_1)$ in Example 3; state s_1 (left) and state s_2 (right)

Then, the functions $g_A(s, f) := h_A(s, f, \gamma(s, f, v), v_A^*)$ for both states s are represented in Figure 14. It is concluded that

$$R_A(s_1, v) = \{(f_1^*, 1 - f_1^*), (1, 0)\} \quad R_A(s_2, v) = \{(f_2^*, 1 - f_2^*)\}.$$

So indeed, the policy f^* of (C.15) is such that $f^*(s) \in R_A(s, v)$ for all s .

There remains to check that the values v_A^* and v_B^* provided are indeed the values of the policy (f^*, g^*) . Indeed:

$$\begin{aligned} h_A(s_1, f^*, b_2, v_A^*) &= \beta_A v_A^*(s_2) = \frac{-3 + \sqrt{19}}{5} = v_A^*(s_1) \\ h_A(s_2, f^*, b_2, v_A^*) &= 1 - (5 - \sqrt{19}) + \frac{1}{2} \left[(5 - \sqrt{19}) \frac{-3 + \sqrt{19}}{5} + (1 - 5 + \sqrt{19}) \frac{-6 + 2\sqrt{19}}{5} \right] \\ &= \sqrt{19} - 4 + \frac{1}{2} \frac{-3 + \sqrt{19}}{5} [-3 + \sqrt{19}] = \sqrt{19} - 4 + \frac{19 + 19 - 6\sqrt{19}}{5} \\ &= -\frac{6}{5} + \frac{2}{5} \sqrt{19} = v_A^*(s_2) \\ h_B(s_1, f, b_2, v_B^*) &= 1 + \beta_B v_B^*(s_2) = 1 + \frac{11 - 2\sqrt{19}}{5} = \frac{16 - 2\sqrt{19}}{5} = v_B^*(s_1) \\ h_B(s_2, f, b_2, v_B^*) &= 2(5 - \sqrt{19}) - 1 + \frac{1}{2} \left[(5 - \sqrt{19}) \frac{16 - 2\sqrt{19}}{5} + (1 - 5 + \sqrt{19}) \frac{22 - 4\sqrt{19}}{5} \right] \\ &= 9 - 2\sqrt{19} + \frac{1}{5} [(5 - \sqrt{19})(8 - \sqrt{19}) - (4 - \sqrt{19})(11 - 2\sqrt{19})] \\ &= 9 - 2\sqrt{19} + \frac{1}{5} [-4 - 19 + 6\sqrt{19}] = \frac{22}{5} - \frac{4}{5} \sqrt{19} = v_B^*(s_2). \end{aligned}$$

We have therefore proved that $v^* = V^{f^* g^*}$, $g^*(s) = \gamma(s, f^*, v^*)$ and $f^* \in R_A(s, v^*)$ for all s . The solution proposed is then a FPE.

References

- [1] Arunabha Bagchi. “Some Economic Applications of Dynamic Stackelberg Games.” In: *Dynamic Games and Applications in Economics*. Ed. by Tamer Başar. Vol. 265. Lecture Notes in Economics and Mathematical Systems. Berlin, Heidelberg: Springer, 1986, pp. 88–102.
- [2] Michèle Breton, Abderrahmane Alj, and Alain Haurie. “Sequential Stackelberg Equilibria in Two-Person Games.” In: *J. Optim. Theory Appl.* 59.1 (1988), pp. 71–97.
- [3] Julio B. Clempner and Alexander S. Poznyak. “Stackelberg security games: Computing the shortest-path equilibrium.” In: *Expert Systems with Applications* 42.8 (2015), pp. 3967–3979.
- [4] Julio Clempner and Alexander Poznyak. “Convergence method, properties and computational complexity for Lyapunov games.” In: *International Journal of Applied Mathematics and Computer Science* 21.2 (2011), pp. 349–361.
- [5] Francesco Maria Delle Fave, Albert Xin Jiang, Zhengyu Yin, Chao Zhang, Milind Tambe, Sarit Kraus, and John P Sullivan. “Game-theoretic patrolling with dynamic execution uncertainty and a case study on a real transit system.” In: *Journal of Artificial Intelligence Research* 50 (2014), pp. 321–367.
- [6] Stephan Dempe. *Foundations of bilevel programming*. Springer Science & Business Media, 2002.
- [7] Eric W. Denardo. “Contraction Mappings in the Theory Underlying Dynamic Programming.” In: *SIAM Review* 9.2 (1967), pp. 165–177.
- [8] Jerzy Filar and Koos Vrieze. *Competitive Markov decision processes*. Springer Science & Business Media, 2012.
- [9] D. Kar, T.H. Nguyen, F. Fang, M. Brown, A. Sinha, M. Tambe, and A.X. Jiang. “Trends and Applications in Stackelberg Security Games.” In: *Handbook of Dynamic Game Theory*. Ed. by T. Başar and G. Zaccour. Springer, 2018. Chap. 28, pp. 1223–1269.
- [10] George Leitman. “On generalized Stackelberg strategies.” In: *J. Optim. Theory Appl.* 26.4 (1978), pp. 637–643.
- [11] Joshua Letchford, Liam MacDermed, Vincent Conitzer, Ronald Parr, and Charles L. Isbell. “Computing Optimal Strategies to Commit to in Stochastic Games.” In: *Twenty-Sixth AAAI Conference on Artificial Intelligence*. 2012.
- [12] Tao Li and Suresh P. Sethi. “A Review of Dynamic Stackelberg Game Models.” In: *Discrete and Continuous Dynamical Systems Series B* 22.1 (2017), pp. 125–159.
- [13] Martin L. Puterman. *Markov decision processes*. Wiley-Interscience, 1994.
- [14] Marwan Simaan and Jose Bejar Cruz Jr. “Additional Aspects on the Stackelberg Strategy in Nonzero-Sum Games.” In: *J. Optim. Theory Appl.* 11.6 (1973), pp. 613–626.

- [15] Marwan Simaan and Jose Bejar Cruz Jr. “On the Stackelberg Strategy in Nonzero-Sum Games.” In: *J. Optim. Theory Appl.* 11.5 (1973), pp. 533–555.
- [16] Yevgeniy Vorobeychik and Satinder Singh. “Computing Stackelberg Equilibria in Discounted Stochastic Games.” In: *Twenty-Sixth AAAI Conference on Artificial Intelligence*. <https://www.aaai.org/ocs/index.php/AAAI/AAAI12/paper/view/4811/5686>. Corrected version retrieved online on Oct. 19, 2018. 2012.
- [17] Yevgeniy Vorobeychik, Bo An, Milind Tambe, and Satinder Singh. “Computing solutions in infinite-horizon discounted adversarial patrolling games.” In: *Proc. 24th International Conference on Automated Planning and Scheduling (ICAPS 2014)(June 2014)*. 2014.
- [18] Ward Whitt. “Representation and Approximation of Noncooperative Sequential Games.” In: *SIAM J. Control and Optimization* 18.1 (1980), pp. 33–48.
- [19] Byung-Wook Wie. “Dynamic Stackelberg equilibrium congestion pricing.” In: *Transportation Research Part C: Emerging Technologies* 15.3 (2007), pp. 154–174.

Contents

1	Introduction	3
1.1	Problem Statement	3
1.2	Related bibliography	4
1.3	Contribution	5
1.4	Notation and definitions	6
2	Operators, Fixed Points and Algorithms	8
2.1	Definition of operators	8
2.2	Properties of operators	10
2.3	Value Iteration Algorithms	12
2.4	Policy Iteration	14
2.5	Mathematical Programming Formulations	17
3	Existence Results for SSSE and FPE.	18
3.1	Single-state results	18
3.2	Myopic Follower Strategies	19
3.3	Zero-Sum Games	22
3.4	Team Games	23
3.5	Acyclic Games	24
4	Numerical Examples.	25
4.1	Experimental setup	25
4.2	Example 1: FPE and SSSE coincide and VI converges	26
4.3	Example 2: FPE and SSSE are different	27
4.4	Example 3: FPE exists but VI does not converge to it	28

5	Application: Surveillance in a graph.	31
5.1	Game description	31
5.2	Computational study	33
6	Conclusions and Further Work	37
A	Notation Summary	39
B	Algorithmic Proof for the Static Case.	40
C	Analysis of the examples	41
C.1	Analysis of Example 1	41
C.1.1	Data	41
C.1.2	Computation of the SSSE	41
C.1.3	One-step Value Iteration from the SSSE	43
C.2	Analysis of Example 2	44
C.2.1	Data	44
C.2.2	Computation of the SSSE	44
C.2.3	The one-step Value Iteration operator	52
C.2.4	One-step Value Iteration from the SSSE	56
C.3	Analysis of Example 3	56
C.3.1	Data	56
C.3.2	Computation of the SSSE	57
C.3.3	Verification of the FPE	58



**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399