



Rolling horizon procedures in Semi-Markov Games: The Discounted Case

Eugenio Della Vecchia, Silvia C. Di Marco, Alain Jean-Marie

► To cite this version:

Eugenio Della Vecchia, Silvia C. Di Marco, Alain Jean-Marie. Rolling horizon procedures in Semi-Markov Games: The Discounted Case. [Research Report] RR-8019, 2012. hal-00720351v1

HAL Id: hal-00720351

<https://inria.hal.science/hal-00720351v1>

Submitted on 25 Jul 2012 (v1), last revised 22 May 2014 (v2)

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.



Rolling horizon procedures in Semi-Markov Games: The Discounted Case

Eugenio Della Vecchia , Silvia Di Marco , Alain Jean-Marie

**RESEARCH
REPORT**

N° 8019

July 2012

Project-Team Maestro



Rolling horizon procedures in Semi-Markov Games: The Discounted Case

Eugenio Della Vecchia ^{*}, Silvia Di Marco [†], Alain Jean-Marie [‡]

Project-Team Maestro

Research Report n° 8019 — July 2012 — 25 pages

Abstract: We study the properties of the rolling horizon and the approximate rolling horizon procedures for the case of two-person zero-sum discounted semi-Markov games with infinite horizon, under several assumptions on the reward function, when the state space is a borelian set and the action spaces are considered compact. Under suitable conditions, we prove that the equilibrium is the unique solution of its dynamic programming equation, and we prove bounds which imply the convergence of the procedures when the horizon length tends to infinity. The approach is based on the formalism for Semi-Markov games developed by Luque-Vásquez in [11], together with extensions of the results of Hernández-Lerma and Lasserre [4] for Markov Decision Processes and Chang and Marcus [2] for Markov Games, both in discrete time. In this way we generalize the results on the rolling horizon and approximate rolling horizon procedures previously obtained for discrete-time problems.

Key-words: Semi-Markov games, Rolling horizon procedures

^{*} CONICET - UNR, Argentina

[†] CONICET - UNR, Argentina

[‡] INRIA and LIRMM, CNRS/Université Montpellier 2, 161 Rue Ada, F-34392 Montpellier, ajm@lirmm.fr.

**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Procédures d'horizon roulant dans les jeux semi-Markoviens: le cas actualisé

Résumé : Nous étudions les propriétés de la procédure de décision à horizon roulant et une approximation de cette procédure, pour le cas de jeux semi-Markoviens à somme nulle avec horizon infini et actualisation, sous différentes hypothèses concernant la fonction de récompense, quand l'espace d'états est un ensemble borélien et les espaces d'actions sont compacts. Sous des hypothèses appropriées, nous montrons que l'équilibre est l'unique solution de l'équation de programmation dynamique associée au jeu, puis nous prouvons des bornes d'erreur impliquant la convergence des procédures quand l'horizon de programmation tend vers l'infini. Notre approche est basée sur le formalisme pour les jeux semi-Markoviens développé par Luque-Vásquez [11], joint à des extensions des résultats de Hernández-Lerma et Lasserre [4] pour les processus de décision Markoviens et Chang et Marcus [2] pour les jeux Markoviens, ces deux derniers travaux étant en temps discret. De cette façon, nous généralisons les résultats sur la procédure à horizon roulant obtenus pour les problèmes en temps discret.

Mots-clés : Jeux semi-Markoviens, procédure à horizon roulant

1 Introduction

In this work we analyze two approximation procedures applied to zero-sum semi-Markov games with the expected total discounted reward as the performance criterion. Specifically, we work with methods derived from the Rolling Horizon procedure.

Semi-Markov Games (**SMG**) generalize Markov Games (**MG**) by allowing the decision maker to choose actions whenever the system state changes, modeling the system evolution in continuous-time and allowing the time spent in a particular state to follow an arbitrary probability distribution. The system state may change several times between decision epochs but only the state at a decision epoch is relevant to the decision maker. Semi-Markov games with discounted reward are analyzed in [11] and later, in [9, 12]. All these papers deal with the characterization of the value function as the fixed point of certain dynamical programming operators. Besides, discrete-time Markov Decision Processes (**MDP**) can be seen as a particular case of **MG** where just one player takes decisions.

On the other hand, Rolling Horizon (**RH**) control is an usual procedure for making decisions in many infinite stage decision problems. It is based on choosing the best most immediate action based on the knowledge of the information of the problem just for a certain number of periods in the future. One design issue of the controller will be then to determine how many periods in the future must be taken into account, in order to make the optimal immediate decision [16]. **RH** strategies are largely used in several areas: we can mention here production control problems, stabilization of control systems, and macro-planning problems. The study of this and other applications can be found in [10].

The objective of the present work is to extend the analysis of the accuracy of the **RH** method and of an Approximate Rolling Horizon (**ARH**) method to Semi-Markov Games with the total discounted reward criterion, when the state space is assumed to be Borel, and the action space compact. As mentioned previously, the **SMG** structure was studied from the theoretical viewpoint by [11] and we borrow some of their results. On the one hand, [4, 5, 2] can be seen as particular cases of our results shown here because the two first investigate the accuracy of the **RH** procedure for discrete-time Markov Decision Processes (**MDP**) (just one decision-maker) with bounded and unbounded rewards functions, respectively, and the last one deals with discrete-time zero-sum Markov games with finite state spaces (two players who play a game whose temporal evolution is discrete). On the other hand, this work also generalizes our previous work [3] where we study the **RH** procedure applied to semi-Markov Decision Processes (**SMDP**) (continuous temporal evolution of the problem where just one decision maker acts). Also as a particular case, all the results obtained here apply to continuous-time **MG** where, up to our knowledge, **RH** and **ARH** have not been applied so far.

This paper is organized as follows. In Section 2, we present the model, the notations and we state the assumptions on the data of the problem. In Section 3 we present the performance criterion and the dynamic programming operator for this case, mentioning the results on the associated optimality equation and recursion scheme. Section 4 contains our contributions about the convergence of **RH** and **ARH** policies rewards to the optimal reward. The approach in this section is based on [5] where the discrete-time **MDP** case is treated and on [2] for the case of finite state spaces. As compared with this last work, we have not only proved results of the **ARH** for more general state spaces and possibly unbounded cost functions, but we also have improved significantly the error bounds. Finally, Section 5 is devoted to the concluding remarks.

2 Preliminaries and Notations

We consider a semi-Markov game of the form

$$M := (\mathcal{S}, \mathcal{A}, \mathcal{B}, \{\mathcal{A}_s : s \in \mathcal{S}\}, \{\mathcal{B}_s : s \in \mathcal{S}\}, Q, F, \ell, \alpha)$$

where

- \mathcal{S} is the state space.
- \mathcal{A} and \mathcal{B} are the action spaces for players 1 and 2 respectively.
- For every $s \in \mathcal{S}$, we define the sets \mathcal{A}_s and \mathcal{B}_s as the sets of actions available in state s for respective players, and, therefore

$$\mathcal{A} = \bigcup_{s \in \mathcal{S}} \mathcal{A}_s, \quad \mathcal{B} = \bigcup_{s \in \mathcal{S}} \mathcal{B}_s.$$

- We set $\mathbb{K} = \{(s, a, b) : s \in \mathcal{S}, a \in \mathcal{A}_s, b \in \mathcal{B}_s\}$.
- The transition law $Q(\cdot|\cdot)$, is a stochastic kernel on \mathcal{S} given \mathbb{K} , and $F(\cdot|s, a, b)$ is the distribution function of the holding time in state $s \in \mathcal{S}$ when actions $a \in \mathcal{A}_s$ and $b \in \mathcal{B}_s$ are chosen.
- The reward function is $\ell : \mathbb{K} \mapsto \mathbb{R}$ and α is a discount factor.

If at time of the n -th decision epoch, the state of the system is $s_n = s$, and the chosen actions are $a_n \in \mathcal{A}_s$ and $b_n \in \mathcal{B}_s$, then the system remains in the state s during a nonnegative random time δ_{n+1} with distribution $F(\cdot|s_n, a_n, b_n)$ and a continuously discounted, stationary reward $\ell(s_n, a_n, b_n)$ is received by player 1 and paid by player 2. The decision epochs are therefore $T_n := T_{n-1} + \delta_n$ for $n \geq 1$, and $T_0 = 0$. The random variable $\delta_{n+1} = T_{n+1} - T_n$ is called the sojourn or holding time at stage n .

For Borel sets X and Y , we will note with $\mathbb{P}(X)$ the family of probability measures on X endowed with the weak topology, and with $\mathbb{P}(X|Y)$ the family of transition probabilities from Y to X .

The space H_n of admissible histories of the process at the n -th decision epoch, consists of sequences of states and decisions up to that epoch. At the initial epoch T_0 , the history consists of the initial state $s_0 \in \mathcal{S}$. At the first decision epoch T_1 , the two initial actions chosen by the players, the holding time at initial state and the new state are added to the initial state, and so on. A typical element of $H_n = (\mathbb{K} \times \mathbb{R}^+)^n \times \mathcal{S}$ is therefore written as

$$h_n = (s_0, a_0, b_0, \delta_1, s_1, a_1, b_1, \delta_2, \dots, s_{n-1}, a_{n-1}, b_{n-1}, \delta_n, s_n).$$

A Markov strategy (or Markov policy) for player 1 is a sequence $\pi = \{\pi_n\}$ of stochastic kernels $\pi_n \in \mathbb{P}(\mathcal{A}|H_n)$ such that for every $h_n \in H_n$ and $n \in \mathbb{N}$, $\pi_n(\mathcal{A}_{s_n}|h_n) = 1$. We denote by Π the set of all Markov strategies of player 1. A Markov strategy $\pi = \{\pi_n\}$ is called stationary if there exists $f \in \mathbb{P}(\mathcal{A}|\mathcal{S})$ such that $f(s) \in \mathbb{P}(\mathcal{A}_s)$ and $\pi_n = f$ for all $s \in \mathcal{S}$ and $n \in \mathbb{N}$. In this case, we identify π with f , i.e., $\pi = f = \{f, f, \dots\}$. We denote by Π_{stat} the set of all stationary strategies.

Similarly, a Markov strategy for player 2 is a sequence $\gamma = \{\gamma_n\}$, where $\gamma_n \in \mathbb{P}(\mathcal{B}|H_n)$, such that for every $h_n \in H_n$ and $n \in \mathbb{N}$, $\gamma_n(\mathcal{B}_{s_n}|h_n) = 1$. In this case we note with Γ the set of all Markov strategies of player 2. A Markov strategy γ is called stationary if there exists $g \in \mathbb{P}(\mathcal{B}|\mathcal{S})$

such that $g(s) \in \mathbb{P}(\mathcal{B}_s)$ and $\gamma_n = g$ for all $s \in \mathcal{S}$ and $n \in \mathbb{N}$. For player 2, we denote Γ_{stat} the set of all stationary strategies.

For each pair of strategies $\pi \in \Pi$ and $\gamma \in \Gamma$, and any initial state s there exist a unique probability measure $P_s^{\pi, \gamma}$ and stochastic processes $\{S_n\}$, $\{A_n\}$, $\{B_n\}$ and $\{\delta_n\}$. S_n , A_n and B_n represent the state and the actions at the n -th decision epoch. $\mathbb{E}_s^{\pi, \gamma}$ denotes the expectation operator with respect to $P_s^{\pi, \gamma}$.

We note

$$\beta(s, a, b) := \int_0^\infty e^{-\alpha t} F(dt|s, a, b) \quad (1)$$

and

$$\vartheta(s, a, b) = \frac{1 - \beta(s, a, b)}{\alpha} . \quad (2)$$

From here on, we make the following abuse of notation: for each $s \in \mathcal{S}$ and given a pair of probability distributions ξ and ζ on \mathcal{A}_s and \mathcal{B}_s respectively, $\int_{\mathcal{A}_s} \int_{\mathcal{B}_s} h(s, a, b) \zeta(db) \xi(da)$ whenever the integral is well defined, will be denoted by $h(s, \xi, \zeta)$. Also, for a function ϕ defined on \mathbb{K} , we note $\phi(s, f, g)$ instead of $\phi(s, f(s), g(s))$, for given stationary policies f and g .

In order to evaluate the performance of policies, we use a total discounted criterion. We assume that the rewards are continuously discounted, with a discount factor α . More precisely let, for $n \geq 1$, $s \in \mathcal{S}$, $\pi \in \Pi$ and $\gamma \in \Gamma$, the expected n -stage α -discounted reward be defined by

$$\begin{aligned} V_n^{\pi, \gamma}(s) &:= \mathbb{E}_s^{\pi, \gamma} \sum_{k=0}^{n-1} \int_{T_k}^{T_{k+1}} e^{-\alpha t} \ell(S_k, A_k, B_k) dt \\ &:= \mathbb{E}_s^{\pi, \gamma} \sum_{k=0}^{n-1} e^{-\alpha T_k} \frac{1 - e^{-\alpha \delta_{k+1}}}{\alpha} \ell(S_k, A_k, B_k) . \end{aligned}$$

The infinite-horizon total expected α -discounted payoff is

$$V^{\pi, \gamma}(s) := \mathbb{E}_s^{\pi, \gamma} \sum_{k=0}^{\infty} e^{-\alpha T_k} \frac{1 - e^{-\alpha \delta_{k+1}}}{\alpha} \ell(S_k, A_k, B_k) . \quad (3)$$

Alternatively, given a pair of policies π and γ , we can write its reward using the variables β and ϑ and obtain for the finite stage horizon and for the infinite horizon respectively,

$$\begin{aligned} V^{\pi, \gamma}(s) &= \mathbb{E}_s \left[\sum_{t=0}^{\infty} \prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \vartheta(S_t, A_t, B_t) \ell(S_t, A_t, B_t) \right] \\ V_n^{\pi, \gamma}(s) &= \mathbb{E}_s \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \vartheta(S_t, A_t, B_t) \ell(S_t, A_t, B_t) \right] , \end{aligned}$$

where we adopt the usual conventions that $\prod_{k=0}^{-1} X_k = 1$ and $\sum_{t=0}^{-1} Y_t = 0$. Briefly, the first expression above comes from the following considerations.

$$V^{\pi, \gamma}(s) = \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{\infty} \int_{T_t}^{T_t + \delta_{t+1}} e^{-\alpha u} du \ell(S_t, A_t, B_t) \right]$$

$$\begin{aligned}
&= \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{\infty} e^{-\alpha T_t} \frac{1 - e^{-\alpha \delta_{t+1}}}{\alpha} \ell(S_t, A_t, B_t) \right] \\
&= \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{\infty} e^{-\alpha (\sum_{k=0}^t \delta_k)} \frac{1 - e^{-\alpha \delta_{t+1}}}{\alpha} \ell(S_t, A_t, B_t) \right] \\
&= \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{\infty} \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) \vartheta(S_t, A_t, B_t) \ell(S_t, A_t, B_t) \right].
\end{aligned}$$

Similarly, the second one can be justified.

At this point, we observe that we can work with an instantaneous one-step reward functions $r: \mathbb{K} \rightarrow \mathbb{R}$ defined by $r(s, a, b) = \vartheta(s, a, b) \ell(s, a, b)$. We obtain the new expressions

$$V^{\pi, \gamma}(s) = \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{\infty} \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t) \right], \quad (4)$$

$$V_n^{\pi, \gamma}(s) = \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{n-1} \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t) \right]. \quad (5)$$

We shall make further assumptions on the distribution probabilities of the holding time and the instantaneous reward function, under which we will work in the future.

Assumption 1.

- (a) The state space \mathcal{S} is a Borel subset of a complete and separable metric space.
- (b) For each $s \in \mathcal{S}$, the sets \mathcal{A}_s and \mathcal{B}_s are compact.
- (c) For each $s \in \mathcal{S}$, and $b \in \mathcal{B}_s$, $r(s, \cdot, b)$ is upper semi-continuous on \mathcal{A}_s .
- (d) For each $s \in \mathcal{S}$, and $a \in \mathcal{A}_s$, $r(s, a, \cdot)$ is lower semi-continuous on \mathcal{B}_s .
- (e) For each $(s, a, b) \in \mathbb{K}$ and each bounded measurable function v on \mathcal{S} , the function $(a, b) \mapsto \int v(y) Q(dy|s, a, b)$ is continuous on $\mathcal{A}_s \times \mathcal{B}_s$.
- (f) The function $(s, a, b) \mapsto \int_0^\infty t F(dt|s, a, b)$ is continuous on \mathbb{K} .

Remark 2.1. Items (a)-(e) in the previous assumption are those considered in [11]. The last one is adopted from [8].

Assumption 2. $\rho := \sup_{(s, a, b) \in \mathbb{K}} \beta(s, a, b) < 1$.

The following result has been extracted from [11, Lemma 4.1].

Proposition 2.1. If there exists a pair of positive numbers θ and ϵ such that

$$F(\theta|s, a, b) \leq 1 - \epsilon$$

for all $(s, a, b) \in \mathbb{K}$, then **Assumption 2** holds with $\rho = 1 - \epsilon + \epsilon e^{\alpha \theta}$.

Let us denote with $\mathcal{M}(\mathcal{S})$ the space of measurable functions on \mathcal{S} , and $\mathcal{M}_+(\mathcal{S})$ the subspace of nonnegative functions of $\mathcal{M}(\mathcal{S})$. On $\mathcal{M}_+(\mathcal{S})$, we define the operator L which maps v to Lv in the following way: for $s \in \mathcal{S}$,

$$(Lv)(s) := \sup_{a \in \mathcal{A}_s, b \in \mathcal{B}_s} \int_{\mathcal{S}} v(z) Q(dz|s, a, b). \quad (6)$$

Let $L^n v = L(L^{n-1}v)$ for $n \in \mathbb{N}$, with $L^0 v = v$. Clearly, $L(\lambda v) = \lambda L(v)$ for every positive scalar λ .

Assumption 3. Let $R(\cdot) : \mathcal{S} \mapsto \mathbb{R}$ be

$$R(s) := \sum_{t=0}^{\infty} \rho^t (L^t r_0)(s) \quad (7)$$

where $r_0(s) = \sup_{a \in \mathcal{A}_s, b \in \mathcal{B}_s} |r(s, a, b)|$. Assume $R(s) < \infty, \forall s \in \mathcal{S}$.

If $\mu \in \mathcal{M}_+(\mathcal{S})$ is strictly positive, for $v \in \mathcal{M}(\mathcal{S})$ we define the μ -weighted norm $\|v\|_\mu = \sup_{s \in \mathcal{S}} |v(s)|/\mu(s)$.

Assumption 4. There exist a measurable function $\mu : \mathcal{S} \rightarrow [1, \infty)$ and a positive constant m such that for all $(s, a, b) \in \mathbb{K}$,

$$(a) \quad |r(s, a, b)| \leq m \mu(s),$$

$$(b) \quad \int \mu(z) Q(dz|s, a, b) \leq \mu(s).$$

Assumption 5. There exists $M > 0$ such that $|r(s, a, b)| \leq M$ for all $(s, a, b) \in \mathbb{K}$.

Remark 2.2. With the definition of r_0 introduced in **Assumption 3**, **Assumption 4 (a)** can be written $\|r_0\|_\mu \leq m$. If r verifies **Assumption 5**, then by setting $\mu \equiv 1$ and $m = M$, **Assumption 4** is satisfied.

In [5] and in [7, Proposition 4.3.1., p. 53], it is shown that **Assumption 4** implies **Assumption 3** for the case of **MDP**. A similar argument is valid for **SMG**'s. Indeed, since

$$|(L^t r_0)(s)| \leq \sup_{a \in \mathcal{A}_s, b \in \mathcal{B}_s} \int_{\mathcal{S}} |(L^{t-1} r_0)(z)| Q(dz|s, a, b),$$

we have

$$\|L^t r_0\|_\mu \leq \|L^{t-1} r_0\|_\mu \leq \dots \leq \|r_0\|_\mu \leq m$$

and then $\|R\|_\mu \leq \frac{m}{1-\rho}$, or, for all $s \in \mathcal{S}$,

$$R(s) \leq \frac{m}{1-\rho} \mu(s)$$

implying the finiteness of R .

Remark 2.3. The theory of **SMDP** can be found in [1] for the case of Borel state spaces, and in [15] and [14] for the case of discrete and finite state spaces, respectively. All of them work with **Assumption 2** and the consequent Proposition 2.1, restricted to only one player.

Besides, an assumption similar to **Assumption 3** appears in [1] for the **SMDP** case and in [5] for the **MDP** case, both for a single player, while **Assumption 4** could be associated to similar ones in [14, Chapter 6] and in [6, Chapter 2], both of them for **MDP**.

Remark 2.4. This semi-Markov environment covers two important special cases:

1. Discrete-time models. In this case $F(\cdot|s, a, b) = \delta_1(\cdot)$ for all $(s, a, b) \in \mathbb{K}$. This corresponds to the theory of **MGs**. Such a function $F(\cdot)$ satisfies **Assumption 1 (f)**.
2. Continuous-time Markov models. This arises if the holding time distributions are exponential: $F(du|s, a, b) = \delta(s, a, b) e^{-\delta(s, a, b)u} du$, where $\delta(s, a, b)$ is a continuous function from \mathbb{K} into $[0, \infty)$ (continuity is required for satisfying **Assumption 1 (f)**). The process $\{S_t\}$ turns out to be a Markov process when (π, γ) is a pair of Markov policies, and a time-homogeneous Markov process when (π, γ) is a pair of stationary policies.

3 Performance Criterion and Related Results

Under suitable hypotheses, in this section we want to characterize the value function for the finite and infinite-horizon games and the strategies which produces the equilibria through operators defined in (10) and (11). We also analyze the convergence of the values for the finite-horizon games to that of the infinite-horizon one. These results will be used to prove the convergence of **RH** and **ARH** procedures described in the next section.

The lower and the upper value functions of the infinite-horizon game are defined, as usual, for $s \in \mathcal{S}$, as

$$\underline{V}^*(s) = \sup_{\pi \in \Pi} \inf_{\gamma \in \Gamma} V^{\pi, \gamma}(s) \quad \text{and} \quad \overline{V}^* = \inf_{\gamma \in \Gamma} \sup_{\pi \in \Pi} V^{\pi, \gamma}(s) \quad (8)$$

respectively. We know that, in general, for all $s \in \mathcal{S}$, $\underline{V}^*(s) \leq \overline{V}^*(s)$. If for all $s \in \mathcal{S}$, $\underline{V}^*(s) = \overline{V}^*(s)$, we refer to this common value as the value of the game, and we note it with V^* . Likewise for the finite-horizon games.

Suppose that our games have a value, then, the objective of the players is to find (when it exists) a pair of policies that solves, given the current state s :

$$(\pi(s), \gamma(s)) = \arg \max_{\pi} \min_{\gamma} V^{\pi, \gamma}(s). \quad (9)$$

Such a pair of strategies $\pi^* \in \Pi$ and $\gamma^* \in \Gamma$ is said to be an equilibrium.

We denote

- $R(\cdot)$ being defined (7), under **Assumption 3**, define $\mathcal{R} = \{v \in \mathcal{M}(\mathcal{S}) : |v(s)| \leq R(s) \text{ for all } s \in \mathcal{S}\}$.
- $\mathcal{M}_{\mu}(\mathcal{S})$ is the linear subspace of $\mathcal{M}(\mathcal{S})$ of the functions with finite μ -weighted norm. This is a Banach space.

Define the operator $T : \mathcal{M}(\mathcal{S}) \mapsto \mathcal{M}(\mathcal{S})$ by

$$(Tv)(s) := \sup_{a \in \mathcal{A}_s} \inf_{b \in \mathcal{B}_s} \left\{ r(s, a, b) + \beta(s, a, b) \int_{\mathcal{S}} v(z) Q(dz|s, a, b) \right\}, \quad (10)$$

and, given a pair of policies $\pi \in \Pi$, $\gamma \in \Gamma$, $T^{\pi, \gamma} : \mathcal{M}(\mathcal{S}) \mapsto \mathcal{M}(\mathcal{S})$

$$(T^{\pi, \gamma}v)(s) := r(s, \pi, \gamma) + \beta(s, \pi, \gamma) \int_{\mathcal{S}} v(z) Q(dz|s, \pi, \gamma). \quad (11)$$

Under **Assumptions 1, 2** and **3**, the sup inf is attained for each $s \in \mathcal{S}$. Denoting with $(f(s), g(s))$ an equilibrium pair at state s , a well-known result of measurable selections (see for example [13, Lemma 4.3]), ascertains that $f, g \in \mathcal{M}(\mathcal{S})$.

Taking into account that $\rho(LR)(s) = \sum_{t=1}^{\infty} \rho^t(L^t r_0)(s) = R(s) - r_0(s)$, note that for all $v \in \mathcal{R}$,

$$|Tv(s)| \leq r_0(s) + \rho(L|v|)(s) \leq r_0(s) + \rho LR(s) = R(s), \quad (12)$$

which implies $Tv \in \mathcal{R}$. On the other hand, under **Assumption 4**, if $v \in \mathcal{M}_{\mu}(\mathcal{S})$, for all $(s, a, b) \in \mathbb{K}$,

$$\int v(z) Q(dz|s, a, b) = \int \frac{v(z)}{\mu(z)} \mu(z) Q(dz|s, a, b) \leq \|v\|_{\mu} \mu(s),$$

and $\|Tv\|_{\mu} \leq m + \rho\|v\|_{\mu}$ which implies that $Tv \in \mathcal{M}_{\mu}(\mathcal{S})$. We have therefore proved the following preliminary result:

Lemma 3.1. *Under **Assumptions 1, 2 and 3**, the operator T maps \mathcal{R} to itself. Under **Assumptions 1, 2 and 4**, T maps $\mathcal{M}_\mu(\mathcal{S})$ to itself. If v is a bounded function, then Tv is also bounded. The same holds for the operators $T^{\pi, \gamma}$.*

The following results provide the existence of optimal stationary strategies, then bounds and convergence results related to the operator L .

Lemma 3.2. *Under **Assumptions 1, 2 and 3**, for each $v \in \mathcal{R}$ there exist stationary strategies $\tilde{f} \in \Pi_{\text{stat}}$, $\tilde{g} \in \Gamma_{\text{stat}}$ such that*

$$\begin{aligned} (Tv)(s) &= r(s, \tilde{f}, \tilde{g}) + \beta(s, \tilde{f}, \tilde{g}) \int_{\mathcal{S}} v(z) Q(dz|s, \tilde{f}, \tilde{g}) \\ &= \sup_{f \in \Pi_{\text{stat}}} \left\{ r(s, f, \tilde{g}) + \beta(s, f, \tilde{g}) \int_{\mathcal{S}} v(z) Q(dz|s, f, \tilde{g}) \right\} \\ &= \inf_{g \in \Gamma_{\text{stat}}} \left\{ r(s, \tilde{f}, g) + \beta(s, \tilde{f}, g) \int_{\mathcal{S}} v(z) Q(dz|s, \tilde{f}, g) \right\}. \end{aligned} \quad (13)$$

Proof. According to Lemma 3.1, for each $v \in \mathcal{R}$, $Tv \in \mathcal{R}$. The proof makes use of well-known measurable selection theorems and is similar to that of Lemma 5.1 in [11], but with the weaker **Assumption 1 (f)**, which is enough to guarantee the existence of the pair of policies stated.

Lemma 3.3. *Suppose that **Assumptions 1, 2 and 3** hold. Then, for all $s \in \mathcal{S}$, $\pi \in \Pi$, $\gamma \in \Gamma$, $v \in \mathcal{R}$ and $t \in \mathbb{N}$, we have*

$$(a) \quad \mathbb{E}_s^{\pi, \gamma}[v(S_t)] \leq (L^t v)(s). \quad (14)$$

$$(b) \quad \lim_{n \rightarrow \infty} \mathbb{E}_s^{\pi, \gamma} \left[\prod_{k=0}^{n-1} \beta(S_k, A_k, B_k) v(S_n) \right] = 0. \quad (15)$$

Proof. We prove (a) by induction on t . Consider the policies $\pi = \{f_0, f_1, \dots\}$ and $\gamma = \{g_0, g_1, \dots\}$. Since we are dealing with Markov policies $\pi \in \Pi$ and $\gamma \in \Gamma$, by [6, Properties 2.5 e and f, pp. 5-6], we have

$$\mathbb{E}_s^{\pi, \gamma}[v(S_{t+1})|S_0 = s, A_0 = a, B_0 = b, S_1 = z] = \mathbb{E}_z^{\pi', \gamma'}[v(S_t)] \quad (16)$$

where we denote with $\pi' = \{f'_0, f'_1, \dots\}$ and $\gamma' = \{g'_0, g'_1, \dots\}$ the policies defined by $f'_t(\cdot|h_t) = f_{t+1}(\cdot|s, a, b, h_t)$ and $g'_t(\cdot|h_t) = g_{t+1}(\cdot|s, a, b, h_t)$.

For $t = 0$ the result is obvious. For $t = 1$, since $\int_{\mathcal{S}} v(z) Q(dz|s, a, b) \leq (Lv)(s)$,

$$\mathbb{E}_s^{\pi, \gamma}[v(S_1)] = \int_{\mathcal{A}} \int_{\mathcal{B}} \int_{\mathcal{S}} v(z) Q(dz|s, a, b) g_0(db) f_0(da) \leq (Lv)(s).$$

In general, using (16), we have for any t :

$$\begin{aligned} \mathbb{E}_s^{\pi, \gamma}[v(S_{t+1})] &= \int_{\mathcal{A}} \int_{\mathcal{B}} \int_{\mathcal{S}} \mathbb{E}_z^{\pi, \gamma}[v(S_{t+1})|S_0 = s, A_0 = a, B_0 = b, S_1 = z] \\ &\quad Q(dz|s, a, b) g_0(db) f_0(da) \end{aligned}$$

$$\begin{aligned}
&= \int_{\mathcal{A}} \int_{\mathcal{B}} \int_{\mathcal{S}} \mathbb{E}_z^{\pi', \gamma'} [v(S_t)] Q(dz|s, a, b) g'_0(db) f'_0(da) \\
&\leq \int_{\mathcal{A}} \int_{\mathcal{B}} \int_{\mathcal{S}} (L^t v)(z) Q(dz|s, a, b) g'_0(db) f'_0(da) \leq (L^{t+1} v)(s).
\end{aligned}$$

For the proof of (b), we have:

$$\begin{aligned}
\mathbb{E}_s^{\pi, \gamma} \left[\left(\prod_{k=0}^{n-1} \beta(S_k, A_k, B_k) \right) v(S_n) \right] &\leq \rho^n \mathbb{E}_s^{\pi, \gamma} v(S_n) \leq \rho^n (L^n v)(s) \\
&\leq \rho^n (L^n R)(s) = \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s) \rightarrow 0,
\end{aligned}$$

as $n \rightarrow \infty$, by **Assumption 3**. \square

The next result generalizes Theorems 4.3 and 4.4 from [11], where proofs are made under **Assumptions 1, 2** and **4**.

Theorem 3.1. *Suppose that **Assumptions 1, 2** and **3** hold. Then*

- (a) *For all $\pi \in \Pi$ and $\gamma \in \Gamma$, $V^{\pi, \gamma}$, \bar{V} and \underline{V} are in \mathcal{R} .*
- (b) *The finite-stage games have values. We have $V_0^* \equiv 0$ and, for $n \geq 1$, the function $V_n^* := TV_{n-1}^* \in \mathcal{R}$ is the value function for the n -stage horizon problem, and the Markovian pair of policies $\{f_0^*, f_1^*, \dots, f_{n-1}^*\}$, $\{g_0^*, g_1^*, \dots, g_{n-1}^*\}$, where the functions f_k^* and g_k^* are the corresponding maxi-minimizing functions, for $k = 0, \dots, n-1$, form an equilibrium.*
- (c) *The infinite-horizon game has a value, and for all $s \in \mathcal{S}$, $|V^*(s) - V_n^*(s)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s) = \rho^n (L^n R)(s) \rightarrow 0$ as $n \rightarrow \infty$.*
- (d) *V^* is the unique function in \mathcal{R} satisfying the optimality equation $TV^* = V^*$. Moreover, there exists a pair of stationary policies (f^*, g^*) which is an equilibrium pair for the infinite-horizon game.*
- (e) *In addition, if **Assumption 4** holds, T is a contractive operator on $\mathcal{M}_\mu(\mathcal{S})$ of modulus ρ and $\|V^* - V_n^*\|_\mu \leq \frac{\rho^n}{1-\rho}$.*

Proof. (a) By definition of r_0 , for all $(s, a, b) \in \mathbb{K}$, $|r(s, a, b)| \leq r_0(s)$, and taking $v = r_0$ in Lemma 3.3 (a), for all $t \in \mathbb{N}$

$$\mathbb{E}_s^{\pi, \gamma} |r(S_t, A_t, B_t)| \leq \mathbb{E}_s^{\pi, \gamma} |r_0(S_t)| \leq (L^t r_0)(s),$$

which implies

$$|V^{\pi, \gamma}(s)| \leq \sum_{t=0}^{\infty} \rho^t (L^t \mathbb{E}_s^{\pi, \gamma} |r(S_t, A_t, B_t)|) \leq R(s)$$

and therefore

$$|\bar{V}(s)| \leq R(s), \quad |\underline{V}(s)| \leq R(s).$$

- (b) The main ideas of the proof can be sketched as follows. By backward induction and well-known arguments, it is possible to prove that starting from $V_0^* \equiv 0 \in \mathcal{R}$, after n successive applications of the operator T , we obtain the value function of the game as well as the equilibrium strategies. In this recursive application it is crucial taking in mind that **Assumption 3** implies that T maps \mathcal{R} into itself (Lemma 3.1). In this way, for each $k \leq n$, V_k^* is in the domain of T , \mathcal{R} .

- (c) For any $s \in \mathcal{S}$, $\pi \in \Pi$ and $\gamma \in \Gamma$, consider formulas (4) and (5) for $V^{\pi, \gamma}(s)$ and $V_n^{\pi, \gamma}(s)$. Then we have

$$|V^{\pi, \gamma}(s) - V_n^{\pi, \gamma}(s)| \leq \sum_{t=n}^{\infty} \rho^t \mathbb{E}_s^{\pi, \gamma} |r(S_t, A_t, B_t)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s). \quad (17)$$

In addition, according to Lemma A.2 in the Appendix,

$$\begin{aligned} |\bar{V}^*(s) - \bar{V}_n^*(s)| &= \left| \inf_{\gamma} \sup_{\pi} V^{\pi, \gamma}(s) - \inf_{\gamma} \sup_{\pi} V_n^{\pi, \gamma}(s) \right| \\ &\leq \sup_{\pi} \sup_{\gamma} |V^{\pi, \gamma}(s) - V_n^{\pi, \gamma}(s)|. \end{aligned}$$

Then using (17),

$$|\bar{V}^*(s) - \bar{V}_n^*(s)| \leq \sup_{\pi} \sup_{\gamma} \sum_{t=n}^{\infty} \rho^t \mathbb{E}_s^{\pi, \gamma} |r(S_t, A_t, B_t)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s).$$

Likewise,

$$|\underline{V}^*(s) - \underline{V}_n^*(s)| \leq \sup_{\pi} \sup_{\gamma} |V^{\pi, \gamma}(s) - V_n^{\pi, \gamma}(s)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s).$$

From part b) we know that each finite-stage game has a value (and then $\underline{V}_n^*(s) = \bar{V}_n^*(s) = V_n^*(s)$), and taking into account the fact that **Assumption 3** implies $\sum_{t=n}^{\infty} \rho^t (L^t r_0)(s) \rightarrow 0$ as $n \rightarrow \infty$, we have that the infinite-horizon game has a value, which verifies

$$V^*(s) = \lim_{n \rightarrow \infty} V_n^*(s),$$

with the convergence bound

$$|V^*(s) - V_n^*(s)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s) = \rho^n (L^n R)(s). \quad (18)$$

- (d) First, we will prove uniqueness. That is, that there is at most one $V \in \mathcal{R}$ such that $TV = V$. For this, observe that Lemma A.2 implies the following inequality for any pair of functions u and v in \mathcal{R} , for any $s \in \mathcal{S}$,

$$|(Tu - Tv)(s)| \leq \rho \sup_{a \in \mathcal{A}_s, b \in \mathcal{B}_s} \int_{\mathcal{S}} |u(z) - v(z)| Q(dz|s, a, b) = \rho (L|u - v|)(s). \quad (19)$$

In general, for $n \in \mathbb{N}$,

$$|(T^n u - T^n v)(s)| \leq \rho^n (L^n |u - v|)(s) \leq 2\rho^n (L^n R)(s) \rightarrow 0,$$

as in Lemma 3.3 (b).

Now we verify the identity proposed. Assume $V^* \in \mathcal{R}$. By Lemma 3.2, there exists a pair of stationary strategies $f^* \in \Pi_{\text{stat}}$ and $g^* \in \Gamma_{\text{stat}}$ such that

$$(TV^*)(s) = r(s, f^*, g^*) + \beta(s, f^*, g^*) \int_{\mathcal{S}} V^*(z) Q(dz|s, f^*, g^*) \quad (20)$$

and V^* is a fixed point in \mathcal{R} of T^{f^*,g^*} .

We show now that for any arbitrary pair $\pi \in \Pi$, $\gamma \in \Gamma$, and $s \in \mathcal{S}$,

$$V^{\pi,g^*}(s) \leq V^{f^*,g^*}(s) \leq V^{f^*,\gamma}(s), \quad (21)$$

which will imply, by (8) and (9), that V^{f^*,g^*} is the value of the game, and that

$$TV^* = T^{f^*,g^*}V^* = T^{f^*,g^*}V^{f^*,g^*} = V^{f^*,g^*} = V^*.$$

The proof of the first inequality in (21) follows passing to the limit as $n \rightarrow \infty$ and from an application of Lemma 3.3 (b) to the following inequality, justified in Lemma 3.4 below:

$$\begin{aligned} V^{f^*,g^*}(s) - \mathbb{E}_s^{\pi,g^*} \left[\left(\prod_{k=0}^n \beta(S_k, A_k, B_k) \right) V^{f^*,g^*}(S_{n+1}) \right] \\ \geq \mathbb{E}_s^{\pi,g^*} \left[\sum_{t=0}^n \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t) \right], \end{aligned}$$

for all $n \in \mathbb{N}$. The second inequality follows in the same manner.

- (e) The contractivity of T is proved in [11, Lemma 6.1] (under a stronger **Assumption 1 f**). Alternatively, in our context, we state the result as a direct consequence of Inequality (19).

For the second part, observe that from (18), $|V^*(s) - V_n^*(s)| \leq \sum_{t=n}^{\infty} \rho^t (L^t r_0)(s)$, and using Remark 2.2

$$\sum_{t=n}^{\infty} \rho^t (L^t r_0)(s) \leq \sum_{t=n}^{\infty} \rho^t m \mu(s) = \frac{m \rho^n}{1 - \rho} \mu(s),$$

which gives

$$\|V^* - V_n^*\|_{\mu} \leq \frac{m \rho^n}{1 - \rho}.$$

Lemma 3.4. *For the pair of strategies $f^* \in \Pi_{\text{stat}}$ and $g^* \in \Gamma_{\text{stat}}$ as in the proof Theorem 3.1 (d), Equation (20), and any strategy $\pi \in \Pi$, there holds:*

$$\begin{aligned} V^{f^*,g^*}(s) - \mathbb{E}_s^{\pi,g^*} \left[\left(\prod_{k=0}^n \beta(S_k, A_k, B_k) \right) V^{f^*,g^*}(S_{n+1}) \right] \\ \geq \mathbb{E}_s^{\pi,g^*} \left[\sum_{t=0}^n \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t) \right]. \end{aligned}$$

Proof. For each $t \geq 1$, $h_t \in H_t$, $a \in \mathcal{A}_{s_t}$, $b \in \mathcal{B}_{s_t}$, since at the t -th epoch of decision, given the history h_t , the realizations of the processes S_k , A_k , and B_k for $k \leq t$ were respectively s_k , a_k and b_k , by properties of the conditional expectation we have

$$\begin{aligned} \mathbb{E}_s^{\pi,g^*} \left[\prod_{k=0}^t \beta(S_k, A_k, B_k) V^{f^*,g^*}(S_{t+1}) | h_t, a_t, b_t \right] \\ = \mathbb{E}_s^{\pi,g^*} \left[\prod_{k=0}^t \beta(s_k, a_k, b_k) V^{f^*,g^*}(S_{t+1}) | h_t, a_t, b_t \right] \end{aligned}$$

$$\begin{aligned}
&= \left(\prod_{k=0}^t \beta(s_k, a_k, b_k) \right) \mathbb{E}_s^{\pi, g^*} \left[V^{f^*, g^*}(S_{t+1}) | h_t, a_t, b_t \right] \\
&= \left(\prod_{k=0}^t \beta(s_k, a_k, b_k) \right) \int_{\mathcal{S}} V^{f^*, g^*}(z) Q(dz | s_t, \pi_n(s_t), g^*(s_t)) \\
&= \left(\prod_{k=0}^{t-1} \beta(s_k, a_k, b_k) \right) \left\{ \beta(s_t, a_t, b_t) \int_{\mathcal{S}} V^{f^*, g^*}(z) Q(dz | s_t, \pi_n(s_t), g^*(s_t)) \right. \\
&\quad \left. + r(s_t, \pi(s_t), g^*(s_t)) - r(s_t, \pi(s_t), g^*(s_t)) \right\} \\
&= \left(\prod_{k=0}^{t-1} \beta(s_k, a_k, b_k) \right) \left\{ (T^{\pi_t(h_t), g^*} V^{f^*, g^*})(s_t) - r(s_t, \pi(s_t), g^*(s_t)) \right\} \\
&\leq \left(\prod_{k=0}^{t-1} \beta(s_k, a_k, b_k) \right) \left\{ V^{f^*, g^*}(s_t) - r(s_t, \pi(s_t), g^*(s_t)) \right\}, \tag{23}
\end{aligned}$$

where, in (22) we make use of the definition of the conditional expectation, and in (23) of Equation (13) from Lemma 3.2. Inequality (23) rewrites as

$$\begin{aligned}
&\left(\prod_{k=0}^{t-1} \beta(s_k, a_k, b_k) \right) V^{f^*, g^*}(s_t) \\
&\quad - \mathbb{E}_s^{\pi, g^*} \left[\left(\prod_{k=0}^t \beta(S_k, A_k, B_k) \right) V^{f^*, g^*}(S_{t+1}) | h_t, a_t, b_t \right] \\
&\geq \left(\prod_{k=0}^{t-1} \beta(s_k, a_k, b_k) \right) r(s_t, \pi(s_t), g^*(s_t)). \tag{24}
\end{aligned}$$

In particular for $t = 0$, this writes as:

$$V^{f^*, g^*}(s_0) - \mathbb{E}_s^{\pi, g^*} \left[\beta(s_0, a_0, b_0) V^{f^*, g^*}(S_1) \right] \geq r(s_0, \pi_0(s_0), g^*(s_0)).$$

Finally, taking expectations under the pair of policies (π, g^*) and summing the last inequality to the ones in (24), written for $t = 1, \dots, n$, we obtain a telescopic sum and therefore the inequality to be proved. \square

4 Approximation Procedures

4.1 Rolling Horizon Procedure

For a wide class of stochastic control or game problems, obtaining an optimal policy explicitly is a difficult task. This is why practitioners often use instead a heuristic method called the Rolling Horizon procedure (also, Receding Horizon, or Model Predictive Control), which works as follows. To the infinite-horizon game is associated a finite-stage horizon¹ game: for a given

¹For continuous-time models, it is important to distinguish a time horizon and an horizon measured in a number of control/game stages.

integer N (the stage-horizon length) and a state s , find:

$$(FHP) \quad \sup_{\pi} \inf_{\gamma} \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{N-1} \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t) \right]. \quad (25)$$

Solving this problem for each initial state results in a sequence of pairs of Markovian policies

$$\pi_N^* = (f_N, f_{N-1}, \dots, f_2, f_1), \quad \gamma_N^* = (g_N, g_{N-1}, \dots, g_2, g_1), \quad (26)$$

where $f_1(s_{N-1})$ is the best action to be applied at stage $t = N - 1$ by the first player, and $g_1(s_{N-1})$ for the second one, when only one stage remains to reach the horizon, f_2 and g_2 are the best decision rules to be applied for the players when two stages remain to get the horizon, at time $t = N - 2$, and so on. In particular, $f_N(s_0)$ and $g_N(s_0)$ are the best decision rules to be applied to the initial state s_0 .

The **RH** method prescribes to repeatedly solve a finite-stage horizon problem, taking the current state as initial state. Then, only the first decision will be applied.

Specifically, the procedure to construct a **RH** policy is the following one. Fix some integer N and consider a denumerable set of epochs.

RH1 At iteration t , and for the current state s_t , solve the N -stage game (FHP), taking s_t as initial state. A pair of actions $f_N(s_t)$, $g_N(s_t)$ is obtained.

RH2 Apply $a_t = f_N(s_t)$, $b_t = g_N(s_t)$.

RH3 Observe the achieved state at time $t + 1$: s_{t+1} .

RH4 Set $t := t + 1$ and $s_t := s_{t+1}$ and go to step 1.

The **RH** procedure does not specify how to compute the values $f_N(s_t)$ and $g_N(s_t)$. Its efficiency is based on the idea that computing the values $f_N(s_t)$ and $g_N(s_t)$ alone is usually much easier than computing the N decision rules in (26). On the other hand, the performance of the resulting sequence of decisions is not the optimal one, although the intuition is that when N is “large enough”, the performance should be close to the optimal. The practical issue is then to choose N so as to obtain a proper compromise between precision and the computational effort needed to obtain $f_N(s_t)$ and $g_N(s_t)$. We address this issue through two formal qualitative and quantitative questions. Let $U_N(s)$ be the expected gain achieved by player 1 with the **RH** procedure with horizon length N , starting in state s :

Q1 Under which conditions on the problem is it true that $\lim_{N \rightarrow \infty} U_N(s) = V^*(s)$?

Q2 Given a state s and $\epsilon > 0$, is it possible to compute N such that $|U_N(s) - V^*(s)| < \epsilon$?

In what follows we prove the convergence of the procedure to the value of the original game. The term “convergence” has to be understood in the sense that when the horizon N goes to infinity, the value obtained with the procedure approaches the value of the game. The preliminary observation, classical for studying Rolling Horizon, is that the procedure **RH** effectively implements, for both players, a stationary Markov policy. Since player 1 will repeatedly play according to the state-feedback function f_N , and player 2 will play g_N , we have $U_N = V^{f_N, g_N}$.

Theorem 4.1. *Suppose that Assumptions 1, 2 and 3 hold. Then, for all $s \in \mathcal{S}$,*

$$0 \leq V^*(s) - U_N(s) \leq 2\rho^N (L^N R)(s).$$

Proof. By definition of an optimal N -stage strategy, for all $s \in \mathcal{S}$ and $N \geq 1$,

$$V_N^*(s) = r(s, f_N, g_N) + \beta(s, f_N, g_N) \int_{\mathcal{S}} V_{N-1}^*(z) Q(dz|s, f_N, g_N), \quad (27)$$

and by (5) and Theorem 1, part b),

$$V_{N-1}^*(s) = \sup_{\pi} \inf_{\gamma} \mathbb{E}_s^{\pi, \gamma} \sum_{t=0}^{N-2} \left(\prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) \right) r(S_t, A_t, B_t).$$

If we add and subtract $\prod_{k=0}^{N-2} \beta(S_k, A_k, B_k) r(S_{N-1}, A_{N-1}, B_{N-1})$:

$$\begin{aligned} V_{N-1}^*(s) = & \sup_{\pi} \inf_{\gamma} \mathbb{E}_s^{\pi, \gamma} \left[\sum_{t=0}^{N-1} \prod_{k=0}^{t-1} \beta(S_k, A_k, B_k) r(S_t, A_t, B_t) \right. \\ & \left. - \prod_{k=0}^{N-2} \beta(S_k, A_k, B_k) r(S_{N-1}, A_{N-1}, B_{N-1}) \right]. \end{aligned}$$

Since for all $s \in \mathcal{S}$, $a \in \mathcal{A}_s$ and $b \in \mathcal{B}_s$, by **Assumption 3**, $-r(s, a, b) \leq r_0(s)$ we have:

$$V_{N-1}^*(s) \leq V_N^*(s) + \sup_{\pi} \sup_{\gamma} \mathbb{E}_s^{\pi, \gamma} \left[\prod_{k=0}^{N-2} \beta(S_k, A_k, B_k) r_0(S_{N-1}) \right].$$

Also, taking into account that $r_0 \geq 0$, by Lemma 3.3, $\mathbb{E}_s^{\pi, \gamma} r_0(S_t) \leq (L^t r_0)(s)$, and

$$V_{N-1}^*(s) \leq V_N^*(s) + \rho^{N-1} (L^{N-1} r_0)(s). \quad (28)$$

By using this inequality in (27) and taking into account that $\beta(s, f_N, g_N) \leq \rho$, we have that:

$$\begin{aligned} V_N^*(s) & \leq r(s, f_N, g_N) + \beta(s, f_N, g_N) \int_{\mathcal{S}} (V_N^*(z) + \rho^{N-1} (L^{N-1} r_0)(z)) Q(dz|s, f_N, g_N) \\ & \leq r(s, f_N, g_N) + \rho^N (L^N r_0)(s) + \beta(s, f_N, g_N) \int_{\mathcal{S}} V_N^*(z) Q(dz|s, f_N, g_N). \end{aligned} \quad (29)$$

Written for $V_N^*(z)$, (29) reads as:

$$V_N^*(z) \leq r(z, f_N, g_N) + \rho^N (L^N r_0)(z) + \beta(z, f_N, g_N) \int_{\mathcal{S}} V_N^*(y) Q(dy|z, f_N, g_N).$$

Substituting inside the integral in (29), we obtain

$$\begin{aligned} V_N^*(s) & \leq r(s, f_N, g_N) + \rho^N (L^N r_0)(s) \\ & \quad + \beta(s, f_N, g_N) \int_{\mathcal{S}} \left[r(z, f_N, g_N) + \rho^N (L^N r_0)(z) \right. \\ & \quad \left. + \beta(z, f_N, g_N) \int_{\mathcal{S}} V_N^*(y) Q(dy|z, f_N, g_N) \right] Q(dz|s, f_N, g_N) \\ & \leq r(s, f_N, g_N) + \beta(s, f_N, g_N) \int_{\mathcal{S}} r(z, f_N, g_N) Q(dz|s, f_N, g_N) \\ & \quad + \rho^N (L^N r_0)(s) + \rho^{N+1} (L^{N+1} r_0)(s) \end{aligned}$$

$$+ \int_{\mathcal{S}} \int_{\mathcal{S}} \beta(s, f_N, g_N) \beta(z, f_N, g_N) V_N^*(y) Q(dy|z, f_N, g_N) Q(dz|s, f_N, g_N) .$$

In general, for every $n \geq 1$,

$$\begin{aligned} V_N^*(s) &\leq \mathbb{E}_s^{f_N, g_N} \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, f_N, g_N) r(S_t, f_N, g_N) \right] \\ &\quad + \sum_{t=N}^{N+n-1} \rho^t(L^t r_0)(s) + \mathbb{E}_s^{f_N, g_N} \left[\prod_{t=0}^{n-1} \beta(S_t, f_N, g_N) V_N^*(S_n) \right] . \end{aligned} \quad (30)$$

Let us analyze the terms of the r.h.s. of (30), as $n \rightarrow \infty$. The first one converges to $U_N(s)$. The second one converges to $\sum_{t=N}^{\infty} \rho^t(L^t r_0)(s)$. On the other hand, since, by Theorem 3.1 part (b), for all $s \in \mathcal{S}$, $|V_N^*(s)| \leq R(s)$, the third term satisfies, by Lemma 3.3

$$\begin{aligned} \left| \mathbb{E}_s^{f_N, g_N} \left[\prod_{t=0}^{n-1} \beta(S_t, f_N, g_N) V_N^*(S_n) \right] \right| &\leq \mathbb{E}_s^{f_N, g_N} \left| \prod_{t=0}^{n-1} \beta(S_t, f_N, g_N) V_N^*(S_n) \right| \\ &\leq \rho^n \mathbb{E}_s^{f_N, g_N} |V_N^*(S_n)| \leq \rho^n \mathbb{E}_s^{f_N, g_N} [R(S_n)] \\ &\leq \rho^n (L^n R)(s) , \end{aligned}$$

which converges to zero, since

$$\rho^n (L^n R)(s) = \sum_{t=n}^{\infty} \rho^t(L^t r_0)(s)$$

and $\sum_{t=0}^{\infty} \rho^t(L^t r_0)(s)$ is supposed to be convergent for all $s \in \mathcal{S}$ in **Assumption 3**. Finally (30) implies

$$V_N^*(s) \leq U_N(s) + \sum_{t=N}^{\infty} \rho^t(L^t r_0)(s), \quad (31)$$

and using Theorem 3.1 (c),

$$\begin{aligned} V^*(s) - U_N(s) &\leq V^*(s) - V_N^*(s) + \sum_{t=N}^{\infty} \rho^t(L^t r_0)(s) \\ &\leq \rho^N (L^N R)(s) + \sum_{t=N}^{\infty} \rho^t(L^t r_0)(s) \\ &= 2\rho^N (L^N R)(s) . \end{aligned}$$

□

Remark 4.1. Theorem 4.1 generalizes to **SMGs** the results in [5, Theorem 4.2] for discrete-time MDP.

In order to improve the bounds obtained, we also can assume the following

Assumption 6. For any $s \in \mathcal{S}$, $V_1^*(s) \geq 0$.

Observe that this is the case, for instance, when the reward function r is positive.

Proposition 4.1. *Suppose that **Assumption 6** holds. Then for $n \in \mathbb{N}$ and for all $s \in \mathcal{S}$, $V_n^*(s) \leq V_{n+1}^*(s)$.*

Proof. Since $V_0^* \equiv 0$, we can write **Assumption 6** as $V_0^*(s) \leq V_1^*(s)$, for all $s \in \mathcal{S}$. If this is the case, for any pair of strategies $f \in \Pi_{\text{stat}}$ and $\gamma \in \Gamma_{\text{stat}}$, it is $T^{f,g}V_0^*(s) \leq T^{f,g}V_1^*(s)$, and by Lemma A.1, $TV_0^*(s) \leq TV_1^*(s)$. The last inequality rewrites as $V_1^*(s) \leq V_2^*(s)$, for all $s \in \mathcal{S}$. The result follows now by induction with similar arguments. \square

Corollary 4.1. *If in Theorem 4.1 **Assumption 6** holds, then, for all $s \in \mathcal{S}$,*

$$0 \leq V^*(s) - U_N(s) \leq \rho^N (L^N R)(s).$$

Proof. If **Assumption 6** holds, by Proposition 4.1 it is possible to improve the bound (28) into: $V_{N-1}^*(s) \leq V_N^*(s)$. Then, instead of Inequality (30), we have:

$$V_N^*(s) \leq \mathbb{E}_s^{f_N, g_N} \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, f_N, g_N) r(S_t, f_N, g_N) \right] + \mathbb{E}_s^{f_N, g_N} \left[\prod_{t=0}^{n-1} \beta(S_t, f_N, g_N) V_N^*(S_n) \right].$$

From here, following with the proof of Theorem 4.1, the new bound follows. \square

Corollary 4.2. *Suppose that **Assumptions 1, 2 and 4** hold. Then, for all $s \in \mathcal{S}$,*

$$0 \leq V^*(s) - U_N(s) \leq \frac{2m\rho^N}{1-\rho} \mu(s),$$

and therefore

$$\|V^* - U_N\|_\mu \leq \frac{2m\rho^N}{1-\rho}.$$

If in addition **Assumption 6** holds, then

$$\|V^* - U_N\|_\mu \leq \frac{m\rho^N}{1-\rho}.$$

Proof. Following the proof of Theorem 4.1 up to Inequality (28) we obtain $V_{N-1}^*(s) \leq V_N^*(s) + \rho^{N-1}(L^{N-1}r_0)(s)$. If **Assumption 4** holds, then (see Remark 2.2), $\|L^{N-1}r_0\|_\mu \leq m$ so that

$$V_{N-1}^*(s) \leq V_N^*(s) + m\rho^{N-1}\mu(s) \tag{32}$$

and (29) can be replaced with

$$V_N^*(s) \leq r(s, f_N, g_N) + m\rho^N\mu(s) + \beta(s, f_N, g_N) \int_{\mathcal{S}} V_N^*(y) Q(dy|s, f_N, g_N).$$

Following the line of proof, we obtain that (30) can be replaced with

$$\begin{aligned} V_N^*(s) \leq & \mathbb{E}_s^{f_N, g_N} \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, f_N, g_N) r(S_t, f_N, g_N) \right] \\ & + \sum_{t=N}^{N+n-1} m\rho^t\mu(s) + \mathbb{E}_s^{f_N, g_N} \left[\prod_{t=0}^{n-1} \beta(S_t, f_N, g_N) V_N^*(S_n) \right]. \end{aligned}$$

As in the proof of Theorem 4.1, the first term converges to $U_N(s)$ and the third one, taking into account the fact that **Assumption 4** implies **Assumption 3**, tends to zero as $n \rightarrow \infty$. The third term converges to $\sum_{t=N}^{\infty} m\rho^t \mu(s) = \frac{m\rho^N}{1-\rho} \mu(s)$, and then the last inequality becomes

$$V_N^*(s) \leq U_N(s) + \frac{m\rho^N}{1-\rho} \mu(s).$$

Finally, from Theorem 3.1 (e),

$$V^*(s) - V_N^*(s) \leq \frac{m\rho^N}{1-\rho} \mu(s), \quad (33)$$

and then

$$V^*(s) - U_N(s) \leq \frac{2m\rho^N}{1-\rho} \mu(s).$$

Now, if we have **Assumption 6**, Inequality (32) can be put in the tighter form $V_{N-1}^*(s) \leq V_N^*(s)$, and continuing with the proof up to Inequality (33), for all $s \in \mathcal{S}$, $V_N^*(s) \leq U_N(s)$. Again the result follows combining Theorem 4.1 with this inequality. \square

When, in particular, we require **Assumption 5** instead of **Assumption 4**, bounds in Corollary 4.2 take the following form:

Corollary 4.3. *Suppose that **Assumptions 1, 2** and **Assumption 5** hold, then for all $s \in \mathcal{S}$*

$$0 \leq V^*(s) - U_N(s) \leq \frac{2M\rho^N}{1-\rho},$$

and therefore,

$$\|V^* - U_N\|_{\infty} \leq \frac{2M\rho^N}{1-\rho}.$$

If in addition **Assumption 6** holds,

$$\|V^* - U_N\|_{\infty} \leq \frac{M\rho^N}{1-\rho}.$$

4.2 Approximate Rolling Horizon Procedure

Suppose now that the players do not have exact information about the problem to be solved at **RH1** in the **RH** procedure, but suppose they know or are able to compute an approximation of that value. We are interested in implementing a procedure where this last approximation is used instead of the value function of the game with finite horizon. We want to estimate the error introduced.

Then, for a function V , supposed to be close in some sense to V_{N-1}^* , choose

$$(\tilde{f}_N(s), \tilde{g}_N(s)) \in \arg \max_{a \in \mathcal{A}_s} \min_{b \in \mathcal{B}_s} \left\{ r(s, a, b) + \beta(s, a, b) \int_{\mathcal{S}} V(z) Q(dz|s, a, b) \right\}.$$

Specifically,

ARH1 Choose some function V a priori near V_{N-1}^* where V_{N-1}^* is the $N-1$ -stage value function.

ARH2 At iteration t , and for the current state s_t , solve

$$\max_{a \in \mathcal{A}_{s_t}} \min_{b \in \mathcal{B}_{s_t}} \left\{ r(s_t, a, b) + \beta(s_t, a, b) \int_{\mathcal{S}} V(z) Q(dz | s_t, a, b) \right\}.$$

A pair of actions $\tilde{f}_N(s_t), \tilde{g}_N(s_t)$ is obtained.

ARH3 Apply $a_t = \tilde{f}_N(s_t), b_t = \tilde{g}_N(s_t)$.

ARH4 Observe the achieved state at time $t + 1$: s_{t+1} .

ARH5 Set $t := t + 1$ and $s_t := s_{t+1}$ and go to step 2.

We will note with \tilde{U}_N the total discounted reward of the pair of stationary policies \tilde{f}_N and \tilde{g}_N . The next result gives answers to questions **Q1** and **Q2** stated in this section for the sequence of successive rewards \tilde{U}_N .

Theorem 4.2. *Suppose that Assumptions 1, 2 and 4 hold. Given $V \in \mathcal{M}_\mu(\mathcal{S})$ a function such that for some $N \geq 1$, $\|V_{N-1}^* - V\|_\mu \leq \varepsilon$, consider a pair of policies $f \in \Pi_{\text{stat}}$ and $g \in \Gamma_{\text{stat}}$ such that $T^{f,g}V = TV$. Then,*

$$\|V^* - V^{f,g}\|_\mu \leq \frac{2m\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho}.$$

If in addition Assumption 6 holds,

$$\|V^* - V^{f,g}\|_\mu \leq \frac{m\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho}.$$

Proof. We shall prove the first statement by bounding the two terms on the r.h.s. of the inequality

$$\|V^* - V^{f,g}\|_\mu \leq \|V^* - TV\|_\mu + \|TV - V^{f,g}\|_\mu. \quad (34)$$

In order to work with the first term, observe that, from Theorem 3.1 part (e),

$$\|TV_{N-1}^* - TV\|_\mu \leq \rho \|V_{N-1}^* - V\|_\mu \leq \rho\varepsilon, \quad (35)$$

and also

$$\|V^* - V_N^*\|_\mu \leq \frac{m\rho^N}{1-\rho}.$$

The previous inequalities gives then

$$\|V^* - TV\|_\mu \leq \|V^* - TV_{N-1}^*\|_\mu + \|TV_{N-1}^* - TV\|_\mu \leq \frac{m\rho^N}{1-\rho} + \rho\varepsilon. \quad (36)$$

Now we work with the second term in the r.h.s. of (34). Observe that, from Inequalities (32) and (35) for all $s \in \mathcal{S}$,

$$V(s) \leq V_{N-1}^*(s) + \varepsilon\mu(s) \leq V_N^*(s) + (\varepsilon + m\rho^{N-1})\mu(s) \leq (TV)(s) + (\rho\varepsilon + \varepsilon + m\rho^{N-1})\mu(s). \quad (37)$$

Then, taking into account Assumption 4, part (b), for all $s \in \mathcal{S}$,

$$(TV)(s) = (T^{f,g})V(s) = r(s, f, g) + \beta(s, f, g) \int_{\mathcal{S}} V(z) Q(dz | s, f, g)$$

$$\begin{aligned}
&\leq r(s, f, g) + \beta(s, f, g) \int_{\mathcal{S}} [(TV)(z) + (\rho\varepsilon + \varepsilon + m\rho^{N-1})\mu(z)] Q(dz|s, f, g) \\
&\leq r(s, f, g) + \beta(s, f, g) \int_{\mathcal{S}} (TV)(z) Q(dz|s, f, g) + (\rho\varepsilon(1 + \rho) + m\rho^N)\mu(s) \\
&= r(s, f, g) \\
&\quad + \beta(s, f, g) \int_{\mathcal{S}} Q(dz|s, f, g) \left(r(z, f, g) + \beta(z, f, g) \int_{\mathcal{S}} V(y) Q(dy|z, f, g) \right) \\
&\quad + (\rho\varepsilon(1 + \rho) + m\rho^N)\mu(s) \\
&= r(s, f, g) \\
&\quad + \beta(s, f, g) \int_{\mathcal{S}} Q(dz|s, f, g) r(z, f, g) \\
&\quad + \beta(s, f, g) \int_{\mathcal{S}} \beta(z, f, g) \int_{\mathcal{S}} V(y) Q(dz|s, f, g) Q(dy|z, f, g) + (\rho\varepsilon(1 + \rho) + m\rho^N)\mu(s) \\
&\leq r(s, f, g) + \beta(s, f, g) \int_{\mathcal{S}} Q(dz|s, f, g) r(z, f, g) \\
&\quad + \beta(s, f, g) \int_{\mathcal{S}} \beta(z, f, g) \int_{\mathcal{S}} (TV)(y) Q(dz|s, f, g) Q(dy|z, f, g) \\
&\quad + \{ [\rho^2\varepsilon(1 + \rho) + \rho\varepsilon(1 + \rho)] + (m\rho^{N+1} + m\rho^N) \} \mu(s).
\end{aligned}$$

In general, for every $n \geq 1$ and $s \in \mathcal{S}$,

$$\begin{aligned}
TV(s) &\leq \mathbb{E}_s^{f,g} \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, f, g) r(S_t, f, g) \right] + \mathbb{E}_s^{f,g} \left[\prod_{t=0}^{n-1} \beta(S_t, f, g) TV(S_n) \right] \\
&\quad + \{ [\rho\varepsilon(1 + \rho) + \dots + \rho^n\varepsilon(1 + \rho)] + (m\rho^N + \dots + m\rho^{N+n-1}) \} \mu(s). \quad (38)
\end{aligned}$$

In order to study the behavior of the terms in the r.h.s. of the last inequality, note that the first term converges to $V^{f,g}$ as $n \rightarrow \infty$, and that the second one converges to zero, since $TV \in \mathcal{M}_\mu(\mathcal{S})$, as in the proof of Corollary 4.2. Now it follows that for all $s \in \mathcal{S}$,

$$TV(s) - V^{f,g}(s) \leq \left[\frac{\rho\varepsilon(1 + \rho)}{1 - \rho} + \frac{m\rho^N}{1 - \rho} \right] \mu(s).$$

Similarly it can be shown that $TV(s) - V^{f,g}(s) \geq - \left[\frac{\rho\varepsilon(1 + \rho)}{1 - \rho} + \frac{m\rho^N}{1 - \rho} \right] \mu(s)$ for all $s \in \mathcal{S}$.

Finally,

$$\begin{aligned}
\|V^* - V^{f,g}\|_\mu &\leq \|V^* - TV\|_\mu + \|TV - V^{f,g}\|_\mu \\
&\leq \frac{m\rho^N}{1 - \rho} + \rho\varepsilon + \frac{\rho\varepsilon(1 + \rho)}{1 - \rho} + \frac{m\rho^N}{1 - \rho} \\
&= \frac{2m\rho^N}{1 - \rho} + \frac{2\rho\varepsilon}{1 - \rho},
\end{aligned}$$

which gives the desired result.

Now, if in addition **Assumption 6** holds, the preceding proof can be done with the inequality $V(s) \leq (TV)(s) + (\rho\varepsilon + \varepsilon)\mu(s)$ instead of Inequality (37), and the bound takes the form

$$\|V^* - V^{f,g}\|_\mu \leq \frac{m\rho^N}{1 - \rho} + \frac{2\rho\varepsilon}{1 - \rho}.$$

□

Remark 4.2. Observe that the result for the case $\varepsilon = 0$ correspond to the N -horizon **RH** approximation, and it was already proved in Corollary 4.2.

Corollary 4.4. Suppose that **Assumptions 1, 2 and 5** hold. Given $V \in \mathcal{M}(\mathcal{S})$ a bounded function such that for some $N \geq 0$, $\|V_{N-1}^* - V\|_\infty \leq \varepsilon$, consider a pair of policies $f \in \Pi_{\text{stat}}$ and $g \in \Gamma_{\text{stat}}$ such that $T^{f,g}V = TV$. Then,

$$\|V^* - V^{f,g}\|_\infty \leq \frac{2M\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho}.$$

If in addition **Assumption 6** holds,

$$\|V^* - V^{f,g}\|_\infty \leq \frac{M\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho}.$$

So far in this work we have considered the situation where, given an approximate value function V , the maximizer plays the policy $f \in \Pi_{\text{stat}}$ such that, for each $s \in \mathcal{S}$, $T^{f,g}V = TV$. Suppose now that this player chooses any $f \in \Pi_{\text{stat}}$ such that, for all $s \in \mathcal{S}$,

$$(T^fV)(s) = (TV)(s),$$

where

$$(T^fV)(s) := \inf_{b \in \mathcal{B}_s} \left\{ r(s, f, b) + \beta(s, f, b) \int_{\mathcal{S}} V(z) Q(dz|s, f, b) \right\}.$$

This situation corresponds to the *worst-case scenario* for player 1. The next result gives us bounds for this second **ARH** framework.

Theorem 4.3. Suppose that **Assumptions 1, 2 and 4** hold. Given $V \in \mathcal{M}_\mu(\mathcal{S})$ a function such that for some $N \geq 1$, $\|V_{N-1}^* - V\|_\mu \leq \varepsilon$, and a policy $f \in \Pi_{\text{stat}}$ such that for all $s \in \mathcal{S}$, $(T^fV)(s) = (TV)(s)$. Then,

$$\|V^* - V^{f,g^*}\|_\mu \leq \frac{2m\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho},$$

where $g^* \in \Gamma_{\text{stat}}$ is an infinite-horizon equilibrium strategy for player 2.

If **Assumption 6** also holds, then

$$\|V^* - V^{f,g^*}\|_\mu \leq \frac{m\rho^N}{1-\rho} + \frac{2\rho\varepsilon}{1-\rho}.$$

Proof. In this case we shall bound the terms on the r.h.s. of the inequality, with $T^fV = TV$,

$$\|V^* - V^{f,g^*}\|_\mu \leq \|V^* - TV\|_\mu + \|TV - V^{f,g^*}\|_\mu.$$

Note that we already have the bound for the first term from Inequality (36): $\|V^* - TV\|_\mu \leq \frac{m\rho^N}{1-\rho} + \rho\varepsilon$. On the other hand, as in the proof of Theorem 4.2,

$$\begin{aligned} V(s) &\leq V_{N-1}^*(s) + \varepsilon\mu(s) \leq V_N^*(s) + (\varepsilon + m\rho^{N-1})\mu(s) \\ &\leq (TV)(s) + (\varepsilon + \rho\varepsilon + m\rho^{N-1})\mu(s) \\ &= (T^fV)(s) + (\varepsilon(1+\rho) + m\rho^{N-1})\mu(s), \end{aligned} \tag{39}$$

and, for every $s \in \mathcal{S}$,

$$\begin{aligned}
(T^f V)(s) &= \inf_{b \in \mathcal{B}_s} \left\{ r(s, f, b) + \beta(s, f, b) \int_{\mathcal{S}} V(z) Q(dz|s, f, b) \right\} \\
&\leq r(s, f, g^*) + \beta(s, f, g^*) \int_{\mathcal{S}} V(z) Q(dz|s, f, g^*) \\
&\leq r(s, f, g^*) + \beta(s, f, g^*) \int_{\mathcal{S}} [(T^f V)(z) + (\varepsilon(1 + \rho) + m\rho^{N-1}) \mu(z)] Q(dz|s, f, g^*) \\
&\leq r(s, f, g^*) + \beta(s, f, g^*) \int_{\mathcal{S}} (T^f V)(z) Q(dz|s, f, g^*) + \rho (\varepsilon(1 + \rho) + m\rho^{N-1}) \mu(s) .
\end{aligned}$$

In general, for every $n \in \mathbb{N}$, and $s \in \mathcal{S}$,

$$\begin{aligned}
(T^f V)(s) &\leq \mathbb{E}_s^{f, g^*} \left[\sum_{t=0}^{n-1} \prod_{k=0}^{t-1} \beta(S_k, f, g^*) r(S_t, f, g^*) \right] + \mathbb{E}_s^{f, g^*} \left[\prod_{t=0}^{n-1} \beta(S_t, f, g^*) T^f V(S_n) \right] \\
&\quad + [(\rho + \dots + \rho^n) \varepsilon(1 + \rho) + (\rho + \dots + \rho^n) m\rho^{N-1}] \mu(s),
\end{aligned}$$

In this case the first term converges to $V^{f, g^*}(s)$ as $n \rightarrow \infty$ and again the second term converges to zero.

Finally,

$$T^f V(s) - V^{f, g^*}(s) \leq \left[\frac{\rho \varepsilon(1 + \rho)}{1 - \rho} + \frac{m\rho^N}{1 - \rho} \right] \mu(s) ,$$

and

$$\begin{aligned}
\|V^* - V^{f, g^*}\|_{\mu} &\leq \frac{m\rho^N}{1 - \rho} + \rho \varepsilon + \frac{\rho \varepsilon(1 + \rho)}{1 - \rho} + \frac{m\rho^N}{1 - \rho} \\
&= \frac{2m\rho^N}{1 - \rho} + \frac{2\rho \varepsilon}{1 - \rho} .
\end{aligned}$$

Again if **Assumption 6** holds, Inequality (39) rewrites $V(s) \leq (T^f V)(s) + \varepsilon(1 + \rho)\mu(s)$ and

$$\|V^* - V^{f, g^*}\|_{\mu} \leq \frac{m\rho^N}{1 - \rho} + \frac{2\rho \varepsilon}{1 - \rho} .$$

□

Corollary 4.5. *Suppose that **Assumptions 1, 2 and 5** hold. Given $V \in \mathcal{M}(\mathcal{S})$ a function such that for some $N \geq 1$, $\|V_{N-1}^* - V\|_{\infty} \leq \varepsilon$, consider a policy $f \in \Pi_{\text{stat}}$ such that for all $s \in \mathcal{S}$, $(T^f V)(s) = (TV)(s)$. Then,*

$$\|V^* - V^{f, g^*}\|_{\mu} \leq \frac{2M\rho^N}{1 - \rho} + \frac{2\rho \varepsilon}{1 - \rho} ,$$

where $g^* \in \Gamma_{\text{stat}}$ is an equilibrium strategy for player 2.

If **Assumption 6** also holds, then

$$\|V^* - V^{f, g^*}\|_{\mu} \leq \frac{M\rho^N}{1 - \rho} + \frac{2\rho \varepsilon}{1 - \rho} .$$

Remark 4.3. Theorems 4.2 and 4.5 generalize to **SMG** our results for **SMDP** presented in [3] and improve results for finite state discrete-time **MG** described in [2, Theorems 6, 7]. The use of finer bounds even allows us to improve their “ $O((1 - \rho)^{-2})$ ” term into a “ $O((1 - \rho)^{-1})$ ” one.

5 Concluding Remarks

Through this work we have dealt with semi-Markov games models with discounted payoff, under different assumptions on the reward function. We have generalized known properties of the equilibria of games for both the finite-horizon and the infinite horizon case.

In addition, we have studied the performance of the rolling horizon procedure and of an approximate rolling horizon procedure. We have proved the convergence of the values related to the rolling horizon procedure to the optimal reward function. We obtain simple pointwise convergence if **Assumption 3** is verified and pointwise geometrical convergence when **Assumption 4** holds. As a particular case, we have obtained uniform geometrical convergence for the case of uniformly bounded rewards functions, i.e. under **Assumption 5**.

Finally we have discussed an approximate rolling horizon procedure, based on the possibility of the controller of not to having perfect prediction of the future needed to take the best immediate action, but approximations of it. Here we have completed the analysis studying the case when the maximizer deals with the possibility of a *worst-case scenario*.

References

- [1] Bhattacharya R., Majumdar, M.; “Controlled semi-Markov models - the discounted case”. *Journal of Statistical Planning and Inference*, 21, 365–381, 1989.
- [2] Chang H.S., Marcus, S.I.; “Two-person zero-sum games: receding horizon approach”. *IEEE Trans. Automatic Control*, 48, 11, 2003, pp. 1951–1961.
- [3] Della Vecchia E., Di Marco S., Jean-Marie A.; “Rolling Horizon Procedure on Controlled Semi-Markov Models. The Discounted Case”. *Annals of the Simposio Argentino en Investigación Operativa* (CD version), ISSN 1850 - 2865. 2011.
- [4] Hernández-Lerma O, Lasserre J.B.; “Error Bounds for Rolling Horizon Policies in Discrete-Time Markov Control Processes”. *IEEE Trans. Automatic Control*, 35, 10, 1990, pp. 1118–1124.
- [5] Hernández-Lerma O, Lasserre J.B.; “Value iteration and rolling plans for Markov control processes with unbounded rewards”. *J. Math. Anal. Appl.*, 177, 1993, pp. 38–55.
- [6] Hernández-Lerma O., *Adaptive Markov Control Processes*. Springer-Verlag 1989.
- [7] Hernández-Lerma O., Lasserre J.B., *Discrete-Time Markov Control Processes*. Springer-Verlag 1996.
- [8] Jaśkiewicz A., Nowak A.; “Approximations of noncooperative semimarkov games”. *J. Optimization, Theory and Applications*, **131**(1), pp. 115–134, 2006.
- [9] Jaśkiewicz A., Nowak A.; “Stochastic Games with Unbounded Payoffs: Applications to Robust Control in Economics”. *Dyn. Games Appl.*, 1, 2011, 253–279.
- [10] Kwon W., Han S.; *Receding Horizon Control. Predictive Control for State Models*. Advanced Textbooks in Control and Signal Processing, Springer, 2005.
- [11] Luque-Vásquez, F.; “Zero-sum semi-Markov game in Borel spaces with discounted payoff”. Research Report, Departamento de Matemáticas, Universidad de Sonora, Mexico. 2002.

- [12] Minjárez-Sosa J., Luque-Vásquez F., “Two Person Zero-Sum Semi-Markov Games with Unknown Holding Times Distribution on One Side: A Discounted Payoff Criterion”, *Appl. Math. Optim.* 57, 2008, 289–305, DOI 10.1007/s00245-007-9016-7.
- [13] Nowak A., “On zero-sum stochastic games with general state space, I”. *Probab. Math. Statist.* 4, 1984, 13–32.
- [14] Puterman L., *Markov Decision Processes*. Wiley and Sons, 2005.
- [15] Ross, S.; *Applied Probability Models with Optimization Applications*, Holden-Day, 1970.
- [16] Sethi T., Sorger G.; “A theory of rolling horizon decision making”. *Ann. Ops. Res.*, 29, 1, 1991, pp. 387–415.

A Appendix

Lemma A.1. *Consider two functions $f, g : X \times Y \rightarrow \mathbb{R}$ such that, for all $x \in X$ and $y \in Y$, $f(x, y) \leq g(x, y)$. Then*

$$\sup_{x \in X} \inf_{y \in Y} f(x, y) \leq \sup_{x \in X} \inf_{y \in Y} g(x, y) .$$

Proof. Consider the functions $F(x) = \inf_{y \in Y} f(x, y)$ and $G(x) = \inf_{y \in Y} g(x, y)$. For any $x \in X$ fixed, and for all $y \in Y$, $F(x) \leq f(x, y) \leq g(x, y)$, and then $F(x) \leq \inf_{y \in Y} g(x, y) = G(x)$. Consequently, $\sup_{x \in X} F(x) \leq \sup_{x \in X} G(x)$, which is the stated inequality. \square

Lemma A.2. *Consider two functions $f, g : X \times Y \rightarrow \mathbb{R}$. Then:*

$$\left| \inf_{y \in Y} \sup_{x \in X} f(x, y) - \inf_{y \in Y} \sup_{x \in X} g(x, y) \right| \leq \sup_{x \in X} \sup_{y \in Y} |f(x, y) - g(x, y)|$$

and

$$\left| \sup_{x \in X} \inf_{y \in Y} f(x, y) - \sup_{x \in X} \inf_{y \in Y} g(x, y) \right| \leq \sup_{x \in X} \sup_{y \in Y} |f(x, y) - g(x, y)| .$$

Proof. Without losing generality, let us suppose $\inf_{y \in Y} \sup_{x \in X} f(x, y) - \inf_{y \in Y} \sup_{x \in X} g(x, y) \geq 0$. If it is not the case, interchange f with g .

Given $\varepsilon > 0$, take $y^* \in Y$ such that

$$\sup_{x \in X} g(x, y^*) - \varepsilon \leq \inf_{y \in Y} \sup_{x \in X} g(x, y) .$$

Then

$$\inf_{y \in Y} \sup_{x \in X} f(x, y) - \inf_{y \in Y} \sup_{x \in X} g(x, y) \leq \sup_{x \in X} f(x, y^*) - \sup_{x \in X} g(x, y^*) + \varepsilon .$$

Now, taking $x^* \in X$ such that

$$\sup_{x \in X} f(x, y^*) \leq f(x^*, y^*) + \varepsilon ,$$

$$\sup_{x \in X} f(x, y^*) - \sup_{x \in X} g(x, y^*) + \varepsilon \leq f(x^*, y^*) - g(x^*, y^*) + 2\varepsilon ,$$

which implies

$$\left| \inf_{y \in Y} \sup_{x \in X} f(x, y) - \inf_{y \in Y} \sup_{x \in X} g(x, y) \right| \leq \sup_{x \in X} \sup_{y \in Y} |f(x, y) - g(x, y)| + 2\varepsilon .$$

Since ε is any arbitrary positive number, the first inequality is proved.

The second inequality follows in the same manner, taking, for $\varepsilon > 0$, $x^* \in X$ such that

$$\sup_{x \in X} \inf_{y \in Y} f(x, y) \leq \inf_{y \in Y} f(x^*, y) + \varepsilon$$

and $y^* \in Y$ such that

$$g(x^*, y^*) - \varepsilon \leq \inf_{y \in Y} g(x^*, y) .$$

Indeed, if $\sup_{x \in X} \inf_{y \in Y} f(x, y) - \sup_{x \in X} \inf_{y \in Y} g(x, y) \geq 0$,

$$\begin{aligned} \sup_{x \in X} \inf_{y \in Y} f(x, y) - \sup_{x \in X} \inf_{y \in Y} g(x, y) &\leq \inf_{y \in Y} f(x^*, y) - \inf_{y \in Y} g(x^*, y) \\ &\leq f(x^*, y^*) - g(x^*, y^*) + 2\varepsilon , \end{aligned}$$

and

$$\left| \sup_{x \in X} \inf_{y \in Y} f(x, y) - \sup_{x \in X} \inf_{y \in Y} g(x, y) \right| \leq \sup_{x \in X} \sup_{y \in Y} |f(x, y) - g(x, y)| + 2\varepsilon .$$

□



**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