WA7 - 10:45

A SAMPLE-PATH APPROACH TO STOCHASTIC GAMES

Melike Baykal-Gursoy Industrial Engineering Department Rutgers University Piscataway, NJ 08854

Abstract

Two-person zero-sum stochastic games with finite state and action spaces are considered. The expected average payoff criterion is used for multichain structures. In the special case that only one player controls the transitions, it is shown that the optimal stationary policies and the value of the game can be obtained from the optimal solutions to a pair of dual linear programs. A decomposition algorithm is given which produces such optimal stationary policies for both players. In the case that both players control the transitions, a generalized game is obtained, the solution of which gives the optimal policies.

1 Introduction

In this paper we investigate the optimal policies for a two-person zero-sum stochastic game (SG) that was first introduced by Shapley [17]. The game is played sequentially. At each epoch, the game is in one of finitely many states and each player observes the current state and chooses one of finitely many actions. The state of the game and the pair of actions determines:(i) a payoff to be made by player II to player I; (ii) the probability distribution over the states of the game, which provides the transition probabilities to the next state of the game. Stochastic games generalize Markov decision processes, in that MDPs may be viewed as SGs in which one of the players has only one action in each state.

An objective function is defined depending on the evaluation of the payoffs. In a game, player I tries to maximize its gain while player II tries to minimize its loss. Shapley [17] considered the game with discounted payoffs and proved that this game has a value and that both players have optimal stationary policies. Hoffman and Karp [10] studied the long-run average payoff criterion and proved that the game has a value if the transition matrix of each pure policy pair is irreducible. If R_m denotes the payoff at epoch $m \in \mathcal{N}_+$, then the long-run average expected payoff to player I is defined as

$$\phi(\mathbf{u},\mathbf{v}) := \liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} E_{\mathbf{u},\mathbf{v}}[R_m],$$

where the expectation is taken with respect to each player's policy. The problem of which games have a value is solved by Mertens and Neyman [14] who show that every stochastic game has a value, i.e.,

$$\sup_{\mathbf{u}} \inf_{\mathbf{v}} \phi(\mathbf{u}, \mathbf{v}) = \inf_{\mathbf{v}} \sup_{\mathbf{u}} \phi(\mathbf{u}, \mathbf{v}).$$

The average expected payoff problem has also been studied by other authors [1], [5], [11]. In general, there do not exist optimal policies for both players [3]. Bewley and Kohlberg [2] give sufficient conditions for the existence of stationary optimal policies for both players. They show that if SG is unichain or if only one player controls the transitions, then there exists optimal stationary policies for both players. Kallenberg [12] study the case when one player controls the transitions and give an algorithm to locate the optimal stationary policies. This case has also been studied by other authors [6],[7], [8], [9], [15].

We consider the expected average criterion (Bewley and Kohl berg define it as limit average criterion [2]),

CH2642-7/89/0000-0180\$1.00© 1989 IEEE

$$\psi(\mathbf{u},\mathbf{v}) := E_{\mathbf{u},\mathbf{v}}[\liminf_{n\to\infty}\frac{1}{n}\sum_{m=1}^{n}R_{m}].$$
 (1)

Bewley and Kohlberg [2] prove that if there exist optimal stationary policies for the long-run average expected payoff criterion, then these policies are optimal for the limit average criterion. In this paper first, we show in the special case that only one player controls the transitions that optimal stationary policies and the value of the game can be obtained from optimal solutions of a pair of dual linear programs. This problem is considered for multichain SGs. The decomposition approach is used in developing the algorithm to locate optimal policies. For the general case that both players control the transitions, we obtain a generalized game. If a solution exists for this game we can locate optimal stationary policies for both players in the unichain case.

Section 2 presents the stochastic game model with some notations and definitions. In section 3 we study the case of one player controlling the transitions and give a decomposition algorithm to locate the optimal policies. In section 4 we discuss the games with both players controlling the transitions.

2 Stochastic Game Model

Let $\{X_m\}$ denote the state process of the two-person zero-sum game, taking values from a finite state space S. After observing the state of the game at epoch m each player chooses an action from a finite set of actions. Let $\{A_m\}$ and $\{B_m\}$ be the sequence of actions taken by player I and player II, respectively. Let A and B denote the set of available actions for player I and II, respectively.

Player II pays player I a payoff $R_m = r(X_m, A_m, B_m)$ at each epoch m. At any epoch if the system is in state $x \in S$, player I choses action $a \in A$ and player II choses action $b \in B$, a payoff of r(x, a, b) is earned. It is assumed that the payoffs are nonnegative and finite. By the time homogeneity assumption the next state of the game depends only on the present state and actions. In particular, when the system is in state $x \in S$ at epoch m and actions $a \in A$ and $b \in B$ are chosen by player I and II, respectively, then the state at epoch m + 1 is $y \in S$ with transition probabilities P_{xaby} , statisfying

$$\sum_{\boldsymbol{y}} P_{\boldsymbol{x} a b \boldsymbol{y}} = 1, \ P_{\boldsymbol{x} a b \boldsymbol{y}} \geq 0, \ \forall \ \boldsymbol{x}, \boldsymbol{y} \in \mathcal{S} \ , \ \boldsymbol{a} \in \mathcal{A} \ , \ \boldsymbol{b} \in \mathcal{B} \ .$$

 P_{xaby} is refered to as the law of motion and assumed to be known to each player.

The probability space that supports the process $\{X_m, A_m, B_m : m \in \mathcal{N}_+\}$ is defined as follows. The underlying sample-space is $\Omega := \{S \times \mathcal{A} \times \mathcal{B}\}^{\infty}$, so that a typical realization $\omega \in \Omega$ is represented by $\omega := (x_1, a_1, b_1, x_2, a_2, b_2, \ldots)$. Let Ω be equipped with the σ -algebra \mathcal{F} generated by the random variables $\{X_m, A_m, B_m : m \in \mathcal{N}_+\}$.

A decision rule \mathbf{u}^m (respectively \mathbf{v}^m) at epoch m for player I (respectively II) is a mapping from $\{S \times \mathcal{A} \times B\}^{m-1} \times S$ to the set of all probability measures on \mathcal{A} (respectively \mathcal{B}). Let $u_a^m(x_1, \ldots, x_m)$ (respectively $v_b^m(x_1, \ldots, x_m)$) denote the conditional probability of choosing action a (respectively b) at epoch m given the past history (x_1, \ldots, x_m) . A policy for player I (respectively II) is denoted by \mathbf{u} (respectively \mathbf{v}).

Let $\xi \in S$ be a fixed initial state known to the players. Policies chosen by each player u and v induce a probability measure $P_{\mathbf{u},\mathbf{v}}$ on (Ω, \mathcal{F}) through the following equations:

$$\mathsf{P}_{\mathbf{u},\mathbf{v}}\{X_1=\xi\}=1,$$

$$P_{\mathbf{u},\mathbf{v}}\{A_m = a | X_1 = x_1, \dots, X_m = x\} = \mathbf{u}_a^m(x_1, \dots, x),$$

$$P_{\mathbf{u},\mathbf{v}}\{B_m = b | X_1 = x_1, \dots, X_m = x\} = \mathbf{v}_b^m(x_1, \dots, x),$$

 $P_{\mathbf{u},\mathbf{v}}\{X_{m+1}=y|X_1=x_1,\ldots,X_m=x,A_m=a,B_m=b\}=P_{xaby}$

Stationary policies f and h for player I and II, respectively, are vectors with components

$$egin{aligned} u^m_a(x_1,\ldots,x) &= f_{xa} \ \ \forall \ m \in \mathcal{N}_+, \ v^m_b(x_1,\ldots,x) &= h_{xb} \ \ \ \forall \ m \in \mathcal{N}_+. \end{aligned}$$

Let C^1 and $C_5^1(C^2, C_5^2)$ denote the set of all policies and stationary policies for player I (II). Under policies $\mathbf{f} \in C_5^1$ and $\mathbf{h} \in C_5^2$ for player I and II respectively, the state process $\{X_m\}$ is a Markov chain with transition probabilities

$$P_{xy}(\mathbf{f},\mathbf{h}) = \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} P_{xaby} f_{xa} h_{xb}.$$

Define the random variables,

$$W_n(x, a, b) := \frac{1}{n} \sum_{m=1}^n \mathbf{1}(X_m = x, A_m = a, B_m = b).$$

Under stationary policies f and h the process $\{Y_m = (X_m, A_m, B_m)\}_{m=1}^{\infty}$ is also a Markov chain, thus,

$$W(x, a, b) = \lim_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} \mathbb{1}(X_m = x, A_m = a, B_m = b)$$

exists $P_{\mathbf{f},\mathbf{h}}$ -almost surely for all $x \in \mathcal{S}$, $a \in \mathcal{A}$ and $b \in \mathcal{B}$.

A stochastic game is said to be unichain if the Markov chain induced by each pure policy pair $(\mathbf{g}_1, \mathbf{g}_2)$ is unichain. Let $P(\mathbf{f}, \mathbf{h})$ denote the transition matrix induced by stationary policies \mathbf{f} and \mathbf{h} for player I and II respectively. If $P(\mathbf{f}, \mathbf{h})$ is unichain, then there exists a unique probability vector $\pi(\mathbf{f}, \mathbf{h}) = \{\pi_x(\mathbf{f}, \mathbf{h}) : x \in S\}$ independent of initial state.

The long-run average reward of player I is given by

$$R = \liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m).$$

If $\{W_n(x, a, b)\}$ converges for all $x \in S$, $a \in A$ and $b \in B$ to random variable $\{W(x, a, b)\}$ $P_{\mathbf{U}, \mathbf{V}}$ -almost surely, then

$$R = \sum_{\boldsymbol{x},\boldsymbol{a},\boldsymbol{b}} r(\boldsymbol{x},\boldsymbol{a},\boldsymbol{b}) W(\boldsymbol{x},\boldsymbol{a},\boldsymbol{b}), \qquad (2)$$

Pu,v-almost surely.

Player I maximizes his expected reward while player II minimizes his expected loss. Hence, the problem is to find policies u[•] and v[•] such that

$$\psi(\mathbf{u},\mathbf{v}^*) \leq \psi(\mathbf{u}^*,\mathbf{v}^*) \leq \psi(\mathbf{u}^*,\mathbf{v}),$$

where $\psi(\mathbf{u}, \mathbf{v})$ is the expected average reward under player I's policy **u** and player II's policy **v**, i.e.,

$$\psi(\mathbf{u},\mathbf{v}) := E_{\mathbf{u},\mathbf{v}}[R].$$

If $\mathbf{u}^*, \mathbf{v}^*$ satisfy the above equation then \mathbf{u}^* and \mathbf{v}^* are called optimal policies for player I and player II, respectively. And $\psi(\mathbf{u}^*, \mathbf{v}^*)$ is called the value of the game. It is straightforward to prove the following proposition.

Proposition 1 Let $\mathbf{f} \in C_s^i$ and $\mathbf{h} \in C_s^i$ and let $\mathcal{R}^1(\mathbf{f}, \mathbf{h}), \ldots, \mathcal{R}^p(\mathbf{f}, \mathbf{h})$ be the recurrent classes induced by $P(\mathbf{f}, \mathbf{h})$. Denote $\{\pi_x^i(\mathbf{f}, \mathbf{h}): x \in \mathcal{R}^i(\mathbf{f}, \mathbf{h})\}$ for the equilibrium probability vector associated with class $i, i = 1, \ldots, p$. Let

$$\tau := \min\{m : m \in \mathcal{N}_+, X_m \in \bigcup_{i=1}^p \mathcal{R}^i(\mathbf{f}, \mathbf{h})\}.$$

Then,

$$\psi(\mathbf{f},\mathbf{h}) = \sum_{i=1}^{\nu} P_{\mathbf{f},\mathbf{h}} \{ X_{\tau} \in \mathcal{R}^{i}(\mathbf{f},\mathbf{h}) \} \sum_{\boldsymbol{x} \in \mathcal{R}^{i}(\mathbf{f},\mathbf{h})} \sum_{\boldsymbol{a} \in \mathcal{A}} \sum_{\boldsymbol{b} \in \mathcal{B}} r(\boldsymbol{x},\boldsymbol{a},\boldsymbol{b}) \\ \pi_{\boldsymbol{x}}^{i}(\mathbf{f},\mathbf{h}) f_{\boldsymbol{x}\boldsymbol{a}} h_{\boldsymbol{x}\boldsymbol{b}}.$$

If the game is unichain, then

$$\psi(\mathbf{f},\mathbf{h}) = \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} r(x,a,b) \pi_x(\mathbf{f},\mathbf{h}) f_{xa} h_{xb}.$$

3 One Player Controlling the Transition Probabilities

We will first investigate a special case, where the transition probabilities are assumed to be independent of the player II's actions. We consider the problem of finding the value of the game and locating the optimal stationary policies for both players. Throughout this section the following assumption holds.

Assumption: The transition probabilities P_{xaby} do not depend on b for all x, y, a, i.e.,

$$P_{\mathbf{u},\mathbf{v}}\{X_{m+1}=y|X_1=x_1,\ldots,X_m=x,A_m=a,B_m=b\}=P_{xay}$$

This condition implies that transitions are influenced by player I's policy u and the components of the transition matrix under a stationary policy are given as

$$P_{xy}(\mathbf{f},\mathbf{h}) = \sum_{a,b} P_{xay} f_{xa} h_{xb} = P_{xy}(\mathbf{f}),$$

which is independent of second player's policy.

Remark: For $\mathbf{f} \in C_s^1$ and $\mathbf{h} \in C_s^2$ and the recurrent classes induced by $P(\mathbf{f}), \mathcal{R}^1(\mathbf{f}), \dots, \mathcal{R}^p(\mathbf{f})$, the expected average payoff is

$$\psi(\mathbf{f},\mathbf{h}) = \sum_{i=1}^{r} P_{\mathbf{f}} \{ X_{\tau} \in \mathcal{R}^{i}(\mathbf{f}) \} \sum_{\boldsymbol{x} \in \mathcal{R}^{i}(\mathbf{f})} \sum_{\boldsymbol{a} \in \mathcal{A}} \sum_{\boldsymbol{b} \in \mathcal{B}} r(\boldsymbol{x},\boldsymbol{a},\boldsymbol{b}) \pi_{\boldsymbol{s}}^{i}(\mathbf{f}) f_{\boldsymbol{s}\boldsymbol{a}} h_{\boldsymbol{x}\boldsymbol{b}}.$$

If the game is unichain, then

$$\psi(\mathbf{f},\mathbf{h}) = \sum_{x \in a \in \mathcal{A}} \sum_{b \in \mathcal{B}} r(x,a,b) \pi_x(\mathbf{f}) f_{xa} h_{xb}. \Box$$

By the above assumption, one can define the strongly communicating class depending only on the law of motion $\{P_{xay}\}$ as follows:

Definition 1 A set of states C is a strongly communicating class if

(i) There exists a stationary policy f such that C is a recurrent class for the associated probability transition matrix P(f),

(ii) C is not a proper subset of some set D which also satisfies (i).

The optimal policies for player I and II shall be constructed using this decomposition approach. The stochastic game is first restricted to each of the strongly communicating classes C^1, \ldots, C^p . Each restricted game corresponds to a pair of dual linear programs. Based on the value of the restricted game an aggregated game is constructed. Since the maximizing player controls the transitions, the aggregated game corresponds to a maximum average reward MDP. An optimal stationary policy for player I is then obtained by combining the optimal policy for the aggregated game with the optimal solutions for the restricted game. The optimal policy for player II is obtained from the restricted game.

To this end, we first define the restricted game.

3.1 The Restricted Game

For each $i = 1, \ldots, s$ and $x \in C^i$, define the set $\mathcal{F}_x = \{a \in \mathcal{A} : P_{xay} = 0 \text{ for all } y \notin C^i\}$ as a subset of \mathcal{A} , player I's actions. Then the restricted game SG-*i* is defined as follows: (*i*) the state space is C^i ; (*ii*) the action space for player I is \mathcal{F}_x for $x \in C^i$, the action space for player II is \mathcal{B} ; (*iii*) P_{xay} and r(x, a, b) are restricted to the state space C^i and action spaces $\mathcal{F}_x, \mathcal{B}$.

For a fixed SG-i and a fixed initial state $\xi \in C^i$, each pair of policies (\mathbf{u}, \mathbf{v}) induces a probability measure $P_{\mathbf{u},\mathbf{v}}^{ii}$ on (Ω, \mathcal{F}) . The corresponding expected average payoff for SG-i is given by

$$\psi^{i}(\mathbf{u},\mathbf{v}) := E^{i}_{\mathbf{u},\mathbf{v}}[\liminf_{n\to\infty}\frac{1}{n}\sum_{m=1}^{n}r(X_{m},A_{m},B_{m})].$$

For each SG-*i*, consider the following pair of dual linear programs.

Program Q_1^i

$$T^{i} := \max \sum_{x \in \mathcal{C}^{i}} \gamma_{x}$$
(3)

s.t.
$$\sum_{\boldsymbol{x}\in \mathcal{C}^{i}}\sum_{\boldsymbol{a}\in\mathcal{F}^{i}} (\delta_{\boldsymbol{x}\boldsymbol{y}} - P_{\boldsymbol{x}\boldsymbol{a}\boldsymbol{y}}) \boldsymbol{z}_{\boldsymbol{x}\boldsymbol{a}} = 0, \quad \forall \ \boldsymbol{y}\in \ \mathcal{C}^{i}$$
(4)

$$\sum_{\mathbf{z}\in \mathcal{C}^{i}}\sum_{\alpha\in\mathcal{F}^{i}}z_{\mathbf{z}\alpha}=1$$
(5)

$$-\sum_{a\in\mathcal{F}^i} r(x,a,b) z_{xa} + \gamma_x \leq 0, \quad \forall \ b\in\mathcal{B} \ , \ x\in \ \mathcal{C}(6)$$

$$z_{za} \ge 0, \ \forall x \in C^i, \ a \in \mathcal{F}^i$$
 (7)

Program Q_2^i

$$U^i := \min \varphi \tag{8}$$

$$s.t.\varphi + \sum_{\mathbf{y}} (\delta_{\mathbf{x}\mathbf{y}} - P_{\mathbf{x}\mathbf{a}\mathbf{y}})t_{\mathbf{y}} - \sum_{b} r(\mathbf{x}, a, b)s_{\mathbf{x}}(b) \ge 0, \quad (9)$$

$$\sum_{b} s_{x}(b) = 1, \quad \forall x \in C^{i}$$
(10)

$$s_x(b) \ge 0, \ \forall \ x \in \ \mathcal{C}^i, \ b \in \mathcal{B}$$
 (11)

Remark: Program Q_1^i maximizes $\sum_{x \in C^i} \sum_{a \in \mathcal{F}_x} r(x, a, b) z_{xa}$.

Theorem 1 Let \mathbf{f}^* denote a stationary optimal policy for player I and let \mathbf{h}^* denote a stationary policy for player II. For each $i = 1, \ldots, s$ and $\mathbf{u} \in C^1$, $\mathbf{v} \in C^2$

$$P_{\mathbf{u},\mathbf{h}} \cdot \{\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) \le T^i | \Phi_i \} = 1,$$
$$P_{\mathbf{f}} \cdot \bigvee_{n \to \infty} \limsup_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) \ge U^i | \Phi_i \} = 1,$$

Furthermore, for a stationary policy h for player II we have

$$P_{\mathbf{f}^*,\mathbf{h}}\{\lim_{n\to\infty}\frac{1}{n}\sum_{m=1}^n r(X_m,A_m,B_m)\geq U^i|\Phi_i\}=1.$$

Proof: Let

$$W_n(x, a, b) := \frac{1}{n} \sum_{m=1}^n \mathbf{1}(X_m = x, A_m = a, B_m = b),$$

$$Z_n(x,a) := \frac{1}{n} \sum_{m=1}^n \mathbf{1}(X_m = x, A_m = a).$$

By compactness properties there is a subsequence $\{N_k(\omega)\}$ along which the limits $W(x, a, b; \omega)$ and $Z(x, a; \omega)$ exist,

$$\lim_{k\to\infty} W_{N_k}(x,a,b) = W(x,a,b) \ge 0, \quad \forall x \in \mathcal{S}, a \in \mathcal{A}, b \in \mathcal{B},$$
$$\lim_{k\to\infty} Z_{N_k}(x,a) = Z(x,a) \ge 0, \quad \forall x \in \mathcal{S}, a \in \mathcal{A}.$$

Let

$$S_{\boldsymbol{x}}(b) = \frac{W(\boldsymbol{x}, \boldsymbol{a}, \boldsymbol{b})}{Z(\boldsymbol{x}, \boldsymbol{a})}.$$

Fix a policy $\mathbf{u} \in C^1$ and $i = \{1, \ldots, s\}$. Assume that player II uses his optimal policy $\mathbf{h}^* \in C_S^2$. Clearly

$$\sum_{b} h^*_{xb} = 1, \ h^*_{xb} \ge 0, \ \forall x \in \mathcal{S} \ , \ b \in \mathcal{B} \ .$$

Let Γ be the set of all sample paths $\omega = (x_1, a_1, b_1, \ldots,)$ that satisfy

(i)
$$a_n \in \mathcal{F}_{x_n}$$
 for all $n \geq N$, for some $N \in \mathcal{N}$,

(ii) $\sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} P_{xay} Z(x, a; \omega) = \sum_{a \in \mathcal{A}} Z(y, a; \omega)$ for all $y \in \mathcal{S}$ (iii) $W(x, a, b; \omega) = h_{xb}^* Z(x, a; \omega)$

By the Strong Law of Large Numbers for Martingale Differences (e.g., see [13]) and the property that $P_{\mathbf{u}}\{A_m \in \mathcal{F}_{X_m} a.a.\} = 1$ for all $\mathbf{u} \in C^1$,

$$P_{\mathbf{u}}(\Gamma) = 1.$$

Since u and h* are independent and

$$P_{\mathbf{h}} \cdot \{S_x(b) = h_{xb}^*\} = 1$$

$$P_{\mathbf{u},\mathbf{h}^{\bullet}}(\Gamma) = 1.$$

We want to show that all sample paths in the intersection of Γ with the set $\Phi_i = \{X_m \in C^i a.a.\}$ satisfy

$$\liminf_{n\to\infty}\frac{1}{n}\sum_{m}^{n}r(X_m,A_m,B_m)\leq T^i,$$

i.e.,

$$\Phi_i \cap \Gamma \subset \{\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^n r(X_m, A_m, B_m) \leq T^i\}.$$

On the set $\Phi_i \cap \Gamma$

$$\begin{array}{rcl} W(x,a,b) &=& 0 & \forall x \notin \mathcal{C}^i, \ a \notin \mathcal{F}_x \\ Z(x,a) &=& 0 & \forall x \notin \mathcal{C}^i, \ a \notin \mathcal{F}_x, \end{array}$$

giving

$$\sum_{x \in C^i} \sum_{a \in \mathcal{F}_x} Z(x, a) = 1,$$
$$\sum_{x \in C^i} \sum_{a \in \mathcal{F}_x} P_{xay} Z(x, a) = \sum_{a \in \mathcal{F}_y} Z(y, a), \quad \forall y \in C^i.$$

Thus, Z(x, a) satisfies primal conditions over the set $\Phi_i \cap \Gamma$. Hence from the primal program (9) we have

$$\sum_{\epsilon} \sum_{\mathcal{C}^i} \sum_{a \in \mathcal{F}_x} \sum_{b \in \mathcal{B}} r(x, a, b) Z(x, a) h_{xb}^* \leq T^i, \qquad (12)$$

on $\Phi_i \cap \Gamma$. On the other hand

$$\liminf_{n\to\infty}\frac{1}{n}\sum_{m=1}^n r(X_m, A_m, B_m) \leq \lim_{k\to\infty}\frac{1}{N_k}\sum_{m=1}^{N_k} r(X_m, A_m, B_m)$$

$$= \sum_{x \in \mathcal{C}^i} \sum_{a \in \mathcal{F}_x} \sum_{b \in \mathcal{B}} r(x, a, b) Z(x, a) h_{xb}^*$$

holds over the set $\Phi_i \cap \Gamma$. When combined with (12) gives the result.

Now fix a policy $\mathbf{v} \in C^2$ for player II and $i = \{1, \ldots, s\}$. Assume that player I uses his optimal policy $\mathbf{f}^* \in C_S^1$. Let Γ denote the set of all sample paths ω that statisfy

(i) $a_n \in \mathcal{F}_{x_n}$ for all $n \geq N$, for some $N \in \mathcal{N}$.

 $\begin{array}{l} (ii) \ \sum_{x \in \mathcal{S}} \ \sum_{a \in \mathcal{A}} \ P_{xay}Z(x,a;\omega) = \sum_{a \in \mathcal{A}} \ Z(y,a;\omega) \text{ for all } y \in \mathcal{S} \ . \\ (iii) \ W(x,a,b;\omega) = \pi_x(\mathbf{f}^*)f_{xa}^*S_x(b;\omega) \end{array}$

 $P_{\mathbf{f}^{*},\mathbf{v}}(\Gamma) = 1.$

Clearly

On the set $\Phi_i \cap \Gamma$

$$W(x, a, b) = 0, \forall x \notin C^{i}, a \notin \mathcal{F}_{x}$$

$$S_{-}(b) = 0, \forall x \notin C^{i}.$$

giving

$$\sum_{x \in \mathcal{B}} S_x(b) = 1 \quad \forall x \in \mathcal{C}^i.$$

Since $z_{\pi a}^* = \pi_{\pi}(\mathbf{f}^*) f_{\pi a}^*$ satisfies the condition (4) and $S_{\pi}(b)$ satisfies dual condition (13) we have from (12),

$$\sum_{x,a,b} r(x,a,b) z_{xa}^* S_x(b) \leq \varphi + \sum_{x,a,y} (\delta_{xy} - P_{xay}) z_{xa}^* t_y$$

= φ .

Thus, dual program minimizes $\sum_{x,a,b} r(x,a,b) z_{xa}^* S_x(b)$, giving

$$\sum_{x \in \mathcal{C}^i} \sum_{a \in \mathcal{F}_x} \sum_{b \in \mathcal{B}} r(x, a, b) z_{xa}^* S_x(b) \ge U^i$$

on the set $\Phi_i \cap \Gamma$. The result is obtained by considering the following equations.

$$\limsup_{n\to\infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) \ge \lim_{k\to\infty} \frac{1}{N_k} \sum_{m=1}^{N_k} r(X_m, A_m, B_m)$$
$$= \sum_{x\in C} \sum_{i \ a\in\mathcal{F}_x} \sum_{b\in\mathcal{B}} r(x, a, b) z_{xa}^* S_x(b).$$

When second player uses stationary policies **h**, then the limit below exists also on the set $\Phi_i \cap \Gamma$

$$\lim_{n\to\infty}\frac{1}{n}\sum_{m=1}^n r(X_m,A_m,B_m) = \sum_{x\in\mathcal{C}^i}\sum_{a\in\mathcal{F}_x}\sum_{b\in\mathcal{B}} r(x,a,b)z_{xa}^*S_x(b).$$

Hence, giving

$$P_{\mathbf{f}^{\bullet},\mathbf{h}}\{\lim_{n\to\infty}\frac{1}{n}\sum_{m=1}^{n}r(X_m,A_m,B_m)\geq U^i\}=1\square$$

Corollary 1 Let \mathbf{f}^* denote the stationary optimal policy for player I, and let \mathbf{h}^* denote the stationary optimal policy for player II. For each SG-i, initial state $\xi \in C^i$, policies $\mathbf{u} \in C^1$ and $\mathbf{v} \in C^2$, $\mathbf{h} \in C^2$

$$P_{\mathbf{u},\mathbf{h}^{*}}^{\ell} \{\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_{m}, A_{m}, B_{m}) \leq T^{i} \} = 1,$$

$$P_{\mathbf{f}^{*},\mathbf{v}}^{\ell} \{\limsup_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_{m}, A_{m}, B_{m}) \geq U^{i} \} = 1,$$

$$P_{\mathbf{f}^{*},\mathbf{h}}^{\ell} \{\lim_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_{m}, A_{m}, B_{m}) \geq U^{i} \} = 1,$$

holds.

Using the definitions of $z_x = \sum_a z_{xa}$ and $I_x^i = \{x \in C^i : z_x > 0\}$, one can construct stationary policies for SG-*i* by an algorithm similar to the one given in [12]

Algorithm 1:

- Find optimal solutions z* = {z*a} and s* = {s*(b)} for the above pair of programs.
- 2. Define $h_{xb}^{*i} := s_x^*(b)$ if $x \in C^i$.
- 3. Define $f_{xa}^{*i} := \frac{x_{xa}}{x_{xa}}$ if $x \in I_x^i$.
- 4. Set $E = I_{*}^{i}$.
- 5. While $E \neq C^i$, Do:
 - Choose a triple (x, a_x, y) such that $x \in C^i E, y \in E$, $a \in \mathcal{F}_x$ and $P_{xa_xy} > 0$.

• Set
$$f_{xa_x}^{\star i} = 1$$
; and Set $E = E \cup \{x\}$.

6. Continue.

Theorem 2 The stationary policies \mathbf{f}^{*i} and \mathbf{h}^{*i} are optimal for SG-i, for all $\xi \in C^i$.

Proof: For a unichain transition matrix

$$z_{xa}^* = \pi_x(\mathbf{f}^{*i}) f_{xa}^{*i}, \quad \forall \ x \in \ \mathcal{C}^i, \ a \in \mathcal{F}_x$$

Also $s^*_x(b) = h^{*i}_{xb}$, for all $x \in C^i$. Since z^*_{xa} and $s^*_x(b)$ are optimal for the pair of mathematical programs, respectively, $T^i = U^i$ holds, giving

$$T^{i} = \sum_{x \in C} \sum_{a \in \mathcal{F}_{x}} \sum_{b \in \mathcal{B}} \tau(x, a, b) z^{*}_{xa} s^{*}_{x}(b)$$
$$= \sum_{x \in C} \sum_{a \in \mathcal{F}_{x}} \sum_{b \in \mathcal{B}} r(x, a, b) \pi_{x}(f^{*i}) f^{*i}_{xa} h^{*i}_{xb}$$
$$= \psi^{i, \ell}(f^{*i}, \mathbf{h}^{*i})$$

By Corollary 1 we have

$$\psi^{i,\ell}(\mathbf{u},\mathbf{h}^{*i}) \leq T^i,$$

$$\psi^{i,\ell}(\mathbf{f}^{*i},\mathbf{h}) > T^i$$

We have

and

$$\psi^{i,\ell}(\mathbf{u},\mathbf{h}^{*i}) \le \psi^{i,\ell}(\mathbf{f}^{*i},\mathbf{h}^{*i}) \le \psi^{i,\ell}(\mathbf{f}^{*i},\mathbf{h}).$$
(13)

Since if one of the players uses a stationary policy, the problem for the other player becomes a MDP and it is well known that in this case there exist a stationary optimal policy, thus implying that it is enough to consider the optimization problem over the stationary policies [12]. Hence (13) implies the result

3.2 The Aggregated SG

Consider the aggregated game, where there is one state corresponding to each strongly communicating class C^i plus states corresponding to the transient states T. For each state $i = 1, \ldots, s$, the action θ is available, which keeps the system in state i with probability 1. Since only the maximizing player controls the transition probabilities, the actions of the form (x, a) are also available, so that the original game moves to state x and player I chooses action $a \notin \mathcal{F}_x$. Thus the aggregated SG is defined so that the state space is $\tilde{S} = \{1, \ldots, s + t\}$ where t denotes the cardinality of the set T, the state-dependent action spaces for player I are

$$\begin{array}{rcl} \mathcal{A}_{\bullet} &=& \{\theta\} \cup \{(x,a): x \in \ \mathcal{C}^{i}, a \notin \mathcal{F}_{x}\} \ i \in \tilde{\mathcal{S}} \ , 1 \leq i \leq s \\ \mathcal{A}_{x} &=& \mathcal{A} \ s+1 \leq i \leq s+t \end{array}$$

the action space for player II is \mathcal{B} , and the law of motion is defined accordingly. The average reward for the aggregated game given that the initial state is $i \in \overline{S}$ is defined as,

$$\bar{\beta}_i(\mathbf{u}) := E_{\mathbf{u}}^i[\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^n \sum_{k=1}^n T^k \ \mathbf{1}(X_m = k)],$$

and the optimal payoff for the aggregated game is

$$\bar{\beta}_i = \sup_{\mathbf{u}\in \mathcal{C}} \bar{\beta}_i(\mathbf{u}).$$

An optimal pure policy g can be found for the aggregated game. The optimal pure policy g of player I moves the game through state space \tilde{S} . Hence, if $g(i) = \theta$ then, it is best to play the restricted game SG-*i* once it has entered the strongly communicating class C^i . But if g(i) = (x, a), then player I moves the original game to state x and then chooses action a.

We can construct the stationary policies as follows. The stationary policy h for player II is obtained by setting

$$h_{xb}^{*} = \begin{cases} h_{xb}^{*i} & \text{if } x \in \mathcal{C}^{i}, \\ 1 & \text{for arbitrary } b \text{ for } x \in \mathcal{T}. \end{cases}$$

And stationary policy f^* for player I can be obtained from the algorithm given in [1]. The following Lemma is stated without proof, since the proof is essentially the same as the proof of Lemma 3.11 in [17]

Lemma 1 The stationary policy f^* constructed by the algorithm is optimal for the intermediate problem. Moreover, if i is a recurrent state under the pure policy g in the aggregated SG, then C^i is closed and contains exactly one recurrent class under $P(f^*)$; if $i \notin H$, then

$$P_{\mathbf{f}^*}(\Phi_i) = 0.$$

Theorem 3 The stationary policies f^* and h^* are optimal for the original problem.

Proof: From Proposition 1

$$E_{\mathbf{f}^*,\mathbf{h}^*}[R] = \sum_{i \in H} P_{\mathbf{f}^*}^{\ell} \{ X_r \in \mathcal{C}^i \} \sum_{x \in \mathcal{C}^i} \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} r(x,a,b)$$

$$\pi_{\mathbf{s}}(\mathbf{f}^{*i}, \mathbf{h}^{*i}) f_{\mathbf{z}a}^{*i} h_{\mathbf{z}b}^{*i} = \sum_{i \in H} P_{\mathbf{f}^*}^{\ell} \{ X_{\tau} \in \mathcal{C}^i \} T^i = \sum_{i=1}^p T^i P_{\mathbf{f}^*}^{\ell} \{ \Phi_i \}.$$

But from Lemma 1 and Theorem 1

$$E_{\mathbf{u},\mathbf{h}^{\bullet}}[R] = \sum_{i=1}^{p} E_{\mathbf{u},\mathbf{h}^{\bullet}}^{\ell} [\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) |\Phi_i] P_{\mathbf{u}}^{\ell} \{\Phi_i\}$$

$$\leq \sum_{i=1}^{p} T^i P_{\mathbf{u}}^{\ell} \{\Phi_i\} \leq E_{\mathbf{f}^{\bullet},\mathbf{h}^{\bullet}}[R].$$

On the other hand,

$$E_{\mathbf{f}^{\bullet},\mathbf{h}}[R] = \sum_{i=1}^{p} E_{\mathbf{f}^{\bullet},\mathbf{h}}^{\ell}[\liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) | \Phi_i] P_{\mathbf{f}^{\bullet}}^{\ell} \{\Phi_i\}$$
$$\geq \sum_{i=1}^{p} T^i P_{\mathbf{f}^{\bullet}}^{\ell} \{\Phi_i\} \square$$

4 General Stochastic Games

In this section we investigate the stochastic game where both players control the transition probabilities. We assume that the Markov chain generated by each stationary policy pair (f,h) is irreducible, i.e., there is only one recurrent class with no transient states.

Assume that player I uses a stationary policy f, then the payoff earned from player II is

$$r_{xb}(\mathbf{f}) := \sum_{a} r(x, a, b) f_{xa},$$

and the transition probabilities determined by player II is

$$P_{xby}(\mathbf{f}) := \sum P_{xaby} f_{xa}.$$

Then, player II will try to minimize his loss, which gives rise to the following mathematical program

$$\begin{split} \min \sum_{x,b} r_{xb}(\mathbf{f}) v_{xb} \\ s.t. \quad \sum_{x,b} P_{xby}(\mathbf{f}) v_{xb} = \sum_{b} v_{yb}, \quad \forall \ y \in \mathcal{S} \\ \sum_{x,b} v_{xb} = 1, \\ v_{xb} \geq 0, \quad \forall \ x \in \mathcal{S} \ , \ b \in \mathcal{B} \ . \end{split}$$

Now let

$$a = \frac{z_{za}}{\sum_{a} z_{za}}$$

f2

with $z_{xa} \ge 0$, and $\sum_a z_{xa} = \sum_b w_{xb}$ for all $x \in S$. And consider the following minimax problem.

$$T := \max_{\mathbf{x}} \min_{\mathbf{v}} \sum_{\mathbf{x}, a, b} r(\mathbf{x}, a, b) \frac{z_{\mathbf{x}a}}{\sum_{a} z_{\mathbf{x}a}} v_{\mathbf{x}b}$$
$$= \min_{\mathbf{v}} \max_{\mathbf{x}} \sum_{\mathbf{x}, a, b} r(\mathbf{x}, a, b) \frac{z_{\mathbf{x}a}}{\sum_{a} z_{\mathbf{x}a}} v_{\mathbf{x}b}$$
(14)

$$\sum_{x,a,b} P_{xaby} \frac{z_{xa}}{\sum_a z_{xa}} v_{xb} = \sum_b v_{xb}, \quad \forall \ y \in \mathcal{S}$$
(15)

$$\sum_{a} z_{xa} = \sum_{b} v_{xb}, \quad \forall x \in \mathcal{S}$$
(16)

$$\sum_{x,b} v_{xb} = 1, \tag{17}$$

$$v_{xb} \geq 0, \ z_{xa} \geq 0, \ \forall x \in S, \ a \in \mathcal{A}, \ b \in \mathcal{B}$$
 (18)

Along the same lines of proof of Theorem 1, one can prove

Theorem 4 Let $\mathbf{f}^* \in C_S^1$ and $\mathbf{h}^* \in C_S^2$ be optimal policies for player I and II, respectively. Then for all policies $\mathbf{f} \in C_S^1$ and $\mathbf{h} \in C_S^2$

$$P_{\mathbf{f}^*,\mathbf{h}}\{\limsup_{n\to\infty}\frac{1}{n}\sum_{m=1}^n r(X_m,A_m,B_m) \ge T\} = 1,$$
$$P_{\mathbf{f},\mathbf{h}^*}\{\liminf_{n\to\infty}\frac{1}{n}\sum_{m=1}^n r(X_m,A_m,B_m) \le T\} = 1.$$

Theorem 5 Suppose that $\{z_{aa}^*\}$ and $\{v_{ab}^*\}$ are solutions to the static game. If f^* and h^* are obtained through the transformations

$$f_{xa}^{*} = \begin{cases} \frac{z^{*}(x,a)}{z^{*}(x)} & \text{if } z^{*}(x) = \sum_{a} z^{*}(x,a) > 0\\ \text{arbitrary but } f_{xa}^{*} > 0, \forall a \in \mathcal{A}, \text{ otherwise} \end{cases}$$
$$h_{xb}^{*} = \begin{cases} \frac{v^{*}(x,b)}{v^{*}(a)} & \text{if } v^{*}(x) = \sum_{b} v^{*}(x,b) > 0\\ \text{arbitrary but } h_{xb}^{*} > 0, \forall b \in \mathcal{B}, \text{ otherwise} \end{cases}$$

then they are optimal for the stochastic game.

Proof: From the constraints we have

$$v^{\bullet}(x) = z^{\bullet}(x),$$

$$\sum_{x} v^{\bullet}(x) = 1,$$

$$v^{\bullet}(y) = \sum_{x,a,b} P_{xaby} v^{\bullet}(x,b) f^{\bullet}_{xa} = \sum_{x,b} P_{xby}(\mathbf{f}^{\bullet}) h^{\bullet}_{xb} v^{\bullet}(x)$$

$$= \sum_{x} P_{xy}(\mathbf{f}^{\bullet},\mathbf{h}^{\bullet}) v^{\bullet}(x).$$

Since there exist a unique probability vector associated with $P(\mathbf{f}^*, \mathbf{h}^*)$, these equations imply that

$$\pi_x(\mathbf{f}^*,\mathbf{h}^*)=v^*(x), \ \forall x\in\mathcal{S}.$$

Thus, we have

$$\begin{split} \liminf_{n \to \infty} \frac{1}{n} \sum_{m=1}^{n} r(X_m, A_m, B_m) &= \sum_{x, a, b} r(x, a, b) \pi_x (\mathbf{f}^*, \mathbf{h}^*) f_{xa}^* h_{xb}^* \\ &= \sum_{x, a, b} r(x, a, b) \frac{v^*(x, b) z^*(x, a)}{v^*(x)} = T \end{split}$$

The result follows from Theorem $1\square$

Thus, we can obtain optimal stationary policies for both players from the static game, if this minimax problem has a solution. **Remark**: Note that in this minimax problem the policies of players are not independent. A game with this added constraint is called *generalized game*[4].

References

- M. Baykal-Gursoy and K. W. Ross. Variance-penalized Markov decision processes. submitted.
- [2] T. Bewley and E. Kohlberg. The asymptotic theory of stochastic games. Mathematics of Operations Research, 1:197-208, 1976.
- [3] T. Bewley and E. Kohlberg. On stochastic games with stationary optimal strategies. Mathematics of Operation Research, 3:104-125, 1978.
- [4] D. Blackwell and T.S. Ferguson. The big match. Annals of Mathematical Statistics, 39:159-163, 1968.
- [5] K.C. Border. Fixed Point Theorems with Applications to Economics and Game Theory. Cambridge Univ. Press, 1985.
- [6] A. Federgruen. successive approximation methods in undiscounted stochastic games. Operations Research, 28:794-809, 1980.
- [7] J.A. Filar. Ordered field property for stochastic games when the player who controls transitions changes from state to state. Journal of Optimization Theory and Appl., 34:503-515, 1981.
- [8] J.A. Filar. On stochastic equilibria of a single-controller stochastic game. Mathematical Programming, 30:313-325, 1984.
- [9] J.A. Filar. Quadratic programming and the single-controller stochastic game. Journal of Math. Analysis and Appl., 113:136-147, 1986.
- [10] J.A. Filar and T.A. Schultz. Nonlinear programming and stationary stategies in stochastic games. *Mathematical Pro*gramming, pages 243-247, 1986.

- [11] A.J. Hoffman and R.M. Karp. On nonterminating stochastic games. Management Science, 12(5):359-370, 1966.
- [12] F. Thuijsman J.A. Filar, T.A. Schultz and D.J. Vrieze. Nonlinear programming and stationary equilibria in stochastic games. Technical report, UMBC, August 1987.
- [13] L.C.M. Kallenberg. Linear Programming and Finite Markovian Control Problems, volume 148. Mathematical Centre Tracts, Amsterdam, 1983.
- [14] M. Loeve. Probability Theory, volume 2. Springer-Verlag, New York, 1978.
- [15] J.F. Mertens and A. Newman. Stochastic games. Int. Journal of Game Theory, 10(2):53-66, 1981.
- [16] T. Parthasarathy and T.E.S. Raghavan. An orderfield property for stochastic games when one player controls transition probabilities. *Journal of Optimization Theory and Applications*, 33(3):375-392, 1981.
- [17] K.W. Ross and R. Varadarajan. The decomposition of time average markov decision processes: Theory, algorithms and applications. *submitted*.
- [18] L.S. Shapley. Stochastic games. Proc. Nat. Acad. Sci. U.S.A., 39:1095-1100, 1953.