

On Markov games with the expected average reward criterion (II)

By

Kensuke TANAKA and Kazuyoshi WAKUTA

(Received September 15, 1975)

1. Introduction

This paper is a continuation of our paper with the same title [2] in which we have showed, mainly, that under some assumptions the Markov games with the criterion of long-run average reward has a value and both players have optimal strategies.

Here, we shall show that the solution $(d, u(\cdot))$ of the functional equation (3.1) of assumption 4 in [2] can be solved by the method of successive approximations under certain conditions. We can obtain this method as an application of the method introduced by D. J. White in Dynamic Programming [3].

2. Preliminaries

In order to state the method of successive approximations, we assume the same conditions as those in [2], that is, (1) S, A and B are compact metric spaces, (2) $r = r(s, a, b)$ is a continuous function on $S \times A \times B$, (3) whenever $s_n \rightarrow s_0, a_n \rightarrow a_0$ and $b_n \rightarrow b_0, q(\cdot | s_n, a_n, b_n)$ converges weakly to $q(\cdot | s_0, a_0, b_0)$, (4) there exist a continuous function $u(s)$ on S and a constant d such that for each $s \in S$,

$$d + u(s) = \sup_{\lambda \in P_A} \inf_{\mu \in P_B} \{r(s, \mu, \lambda) + \int u(s') dq(s' | s, \mu, \lambda)\}, \quad (2.1)$$

where P_A and P_B are the sets of all probability measures on $(A, \mathfrak{B}(A))$ and $(B, \mathfrak{B}(B))$, respectively,

$$r(s, \mu, \lambda) = \iint r(s, a, b) d\mu(a) d\lambda(b),$$

and

$$q(E | s, \mu, \lambda) = \iint q(E | s, a, b) d\mu(a) d\lambda(b).$$

Moreover, we need the same lemmas as those in [2], that is,

LEMMA 1. Let $u(s, a, b)$ be a continuous, real-valued function on $S \times A \times B$. Then $u(s, a, b) = \int_A \int_B u(s, a, b) d\mu(a) d\lambda(b)$, $s \in S$, $\mu \in P_A$, $\lambda \in P_B$, is a continuous function on $S \times P_A \times P_B$.

LEMMA 2. Let $u(x, y)$ be a bounded, continuous function on $X \times Y$, where X is a Borel subset of a Polish space and Y is a compact metric space. Then, $u^*: X \rightarrow R$ defined by $u^*(x) = \max_{y \in Y} u(x, y)$ is continuous. Moreover, $u_*: X \rightarrow R$ defined by $u_*(x) = \min_{y \in Y} u(x, y)$ is also continuous.

LEMMA 3. Let $u(x, y)$ be a bounded, continuous function on $X \times Y$, where X is a Borel subset of a Polish space and Y is a compact metric space. Then, there exist Borel measurable maps f and g from X to Y such that $u(x, f(x)) = \max_{y \in Y} u(x, y)$, $x \in X$ and $u(x, g(x)) = \min_{y \in Y} u(x, y)$, $x \in X$.

The reader is referred to [1] for the proofs,

Then, from Lemma 1 and Lemma 2, the equation (1, 1) can be replaced by the following:

$$\begin{aligned} d+u(s) &= \max_{\mu \in P_A} \min_{\lambda \in P_B} \{r(s, \mu, \lambda) + \int u(s') dq(s' | s, \mu, \lambda)\} \\ &= \min_{\lambda \in P_B} \max_{\mu \in P_A} \{r(s, \mu, \lambda) + \int u(s') dq(s' | s, \mu, \lambda)\}. \end{aligned} \quad (2.2)$$

We want to show that the solution $(u(\cdot), d)$ of (2.2) can be solved by the method of successive approximations.

3. The method of successive approximations

Let Π^n be the set of all Markov policies for player I until the n th time and Γ^n the set of all Markov policies for player II until the n th time. That is, for $\pi \in \Pi^n$ and $\sigma \in \Gamma^n$, $\pi = (f_1, f_2, \dots, f_n)$ and $\sigma = (g_1, \dots, g_n)$, where f_k and g_k , $k=1, 2, \dots, n$, are Borel measurable maps from S to P_A and P_B , respectively. Then, for each $\pi = (f_1, f_2, \dots, f_n) \in \Pi^n$, $\sigma = (g_1, g_2, \dots, g_n) \in \Gamma^n$, $E \in \mathfrak{B}(S)$ and $s_1 \in S$, we define the following:

$$\begin{aligned} q^{(n)}(E | s_1, \pi, \sigma) &= \int \dots \int q(E | s_n, f_n(s_n), g_n(s_n)) \\ &\quad \prod_{i=1}^{n-1} dq(s_{i+1} | s_i, f_i(s_i), g_i(s_i)). \end{aligned} \quad (3.1)$$

And, moreover, we assume the following assumption.

ASSUMPTION. There exist an integer $u \geq 0$, a quantity $\alpha (0 < \alpha \leq 1)$ and a state $s_0 \in S$ such that, for each $\pi = (f_1, f_2, \dots, f_{u+1}) \in \Pi^{u+1}$, $\sigma = (g_1, g_2, \dots, g_{u+1}) \in \Gamma^{u+1}$ and $s_1 \in S$,

$$q^{(u+1)}(s_0 | s_1, \pi, \sigma) \geq \alpha > 0. \quad (3.2)$$

Then, it should be noted that assumption 5 in [2] implies this assumption.

THEOREM. Under the above assumption, the sequence $\{d_n, v_n(\cdot), n \geq 0\}$ defined by, $s \in S$

$$V_n(s) = \max_{\mu \in P_A} \min_{\lambda \in P_B} \{r(s, \mu, \lambda) + \int_S v_{n-1}(s') dq(s' | s, \mu, \lambda)\},$$

$$d_n = V_n(s_0), \tag{3.3}$$

$$v_n(s) = V_n(s) - d_n,$$

converges uniformly to the solution $\{d, u(\cdot)\}$ of the functional equation (2.1), where

$$v_0(s) = \max_{\mu \in P_A} \min_{\lambda \in P_B} r(s, \mu, \lambda).$$

PROOF. By Lemma 1 and Lemma 2, it follows that $V_n(\cdot), v_n(\cdot)$ are well defined and continuous. We need only prove the uniformity of convergence, since this implies that the limiting form is a solution of (2.1). For any sequence of continuous functions $\{Z_n(s)\}$ define

$$\nabla_n(Z) = \inf_{s \in S} [Z_n(s) - Z_{n-1}(s)] \tag{3.4}$$

and

$$\Delta_n(Z) = \sup_{s \in S} [Z_n(s) - Z_{n-1}(s)]. \tag{3.5}$$

When $n \geq u + 4$, we can show that

$$V_n(s) - V_{n-1}(s) \tag{3.6}$$

$$\geq \min_{\mu} [\min_{\lambda} \{r(s, \mu, \lambda) + \int v_{n-1}(s') dq(s' | s, \mu, \lambda)\} - \min_{\lambda} \{r(s, \mu, \lambda) + \int v_{n-2}(s') dq(s' | s, \mu, \lambda)\}]$$

$$\geq \min_{\mu} [\min_{\lambda} \int (v_{n-1}(s') - v_{n-2}(s')) dq(s' | s, \mu, \lambda)]$$

$$= \min_{\mu, \lambda} [\int (v_{n-1}(s') - v_{n-2}(s')) dq(s' | s, \mu, \lambda)]$$

$$= \min_{\mu, \lambda} [\int (V_{n-1}(s') - V_{n-2}(s')) dq(s' | s, \mu, \lambda)] - (d_{n-1} - d_{n-2}).$$

Hence,

$$V_n(s) - V_{n-1}(s) + (d_{n-1} - d_{n-2})$$

$$\geq \min_{\mu, \lambda} [\int (V_{n-1}(s') - V_{n-2}(s')) dq(s' | s, \mu, \lambda)]. \tag{3.7}$$

Then, by Lemma 3, we can show that

$$\begin{aligned}
V_n(s) - V_{n-1}(s) + (d_{n-1} - d_{n-2}) + (d_{n-2} - d_{n-3}) & \quad (3.8) \\
& \geq \min_{\mu, \lambda} \left[\int \left\{ \min_{\mu', \lambda'} \int (V_{n-2}(s'') - V_{n-3}(s'')) dq(s'' | s', \mu', \lambda') \right\} \right. \\
& \quad \left. dq(s' | s, \mu, \lambda) \right]. \\
& \geq \min_{\pi \in \Pi^2, \sigma \in \Gamma^2} \left[\int (V_{n-2}(s') - V_{n-3}(s')) dq^{(2)}(s' | s, \pi, \sigma) \right].
\end{aligned}$$

Repeating this process we derive

$$\begin{aligned}
V_n(s) - V_{n-1}(s) + \sum_{k=1}^u (d_{n-k} - d_{n-k-1}) & \quad (3.9) \\
& \geq \min_{\pi \in \Pi^u, \sigma \in \Gamma^u} \left[\int (V_{n-u}(s') - V_{n-u-1}(s')) dq^{(u)}(s' | s, \pi, \sigma) \right],
\end{aligned}$$

namely,

$$\begin{aligned}
V_n(s) - V_{n-1}(s) + (d_{n-1} - d_{n-u}) & \quad (3.10) \\
& \geq \min_{\pi \in \Pi^{u+1}, \sigma \in \Gamma^{u+1}} \left[\int (v_{n-u-1}(s') - v_{n-u-2}(s')) dq^{(u+1)}(s' | s, \pi, \sigma) \right].
\end{aligned}$$

Now by hypothesis (3.3)

$$v_{n-u-1}(s_0) = v_{n-u-2}(s_0) = 0$$

and

$$q^{(u+1)}(s_0 | s, \pi, \sigma) \geq \alpha > 0 \quad \text{for all } \pi \in \Pi^{u+1}, \sigma \in \Gamma^{u+1} \text{ and } s \in S.$$

Thus, we have

$$\nabla_n(V) + (d_{n-1} - d_{n-u}) \geq (1 - \alpha) \nabla_{n-u-1}(v). \quad (3.11)$$

Similarly, for $n \geq u + 4$,

$$\begin{aligned}
V_n(s) - V_{n-1}(s) + (d_{n-1} - d_{n-u}) & \\
& \leq \max_{\pi \in \Pi^{u+1}, \sigma \in \Gamma^{u+1}} \left[\int (v_{n-u-1}(s') - v_{n-u-2}(s')) \right. \\
& \quad \left. dq^{(u+1)}(s' | s, \pi, \sigma) \right]. & \quad (3.12)
\end{aligned}$$

Thus, by hypothesis (3.3),

$$\Delta_n(V) + (d_{n-1} - d_{n-u}) \leq (1 - \alpha) \Delta_{n-u-1}(v). \quad (3.13)$$

(3.11) and (3.13) imply that

$$\Delta_n(V) - \nabla_n(V) \leq (1 - \alpha) (\Delta_{n-u-1}(v) - \nabla_{n-u-1}(v)). \quad (3.14)$$

At the same time,

$$\nabla_n(v) = \nabla_n(V) - (d_n - d_{n-1})$$

and

$$\Delta_n(v) = \Delta_n(V) - (d_n - d_{n-1})$$

yield

$$\Delta_n(v) - \nabla_n(v) = \Delta_n(V) - \nabla_n(V). \quad (3.15)$$

Thus, by (3.15), (3.14) becomes

$$\Delta_n(v) - \nabla_n(v) \leq (1-\alpha)(\Delta_{n-u-1}(v) - \nabla_{n-u-1}(v)).$$

Let $D_n(v) = \Delta_n(v) - \nabla_n(v)$. Then, for $n = N(u+1) + r$ ($1 \leq r \leq u+1$),

$$D_n(v) \leq (1-\alpha)^N D_r(v) \leq (1-\alpha)^N A,$$

where $A = \max \{D_1(v), D_2(v), \dots, D_{u+1}(v)\}$.

Since $v_n(s_0) - v_{n-1}(s_0) = 0$, it follows that

$$\nabla_n(v) \leq 0 \leq \Delta_n(v).$$

Let

$$U_n(v) = \sup_{s \in S} |v_n(s) - v_{n-1}(s)|$$

and

$$U_n(V) = \sup_{s \in S} |V_n(s) - V_{n-1}(s)|,$$

then

$$U_n(v) \leq D_n(v) \leq (1-\alpha)^N A. \quad (3.16)$$

This is sufficient to prove that the sequence $\{v_n(s)\}$ converges uniformly to the function $v(s)$.

On the other hand, we easily see that

$$U_n(V) \leq U_{n-1}(v). \quad (3.17)$$

From (3.17), the sequence $\{V_n(s)\}$ also converges uniformly to a function $V(s)$. Since $\{V_n(s)\}$ converges uniformly, so does the sequence $\{d_n\} = \{V_n(s_0)\}$ converge. Thus, the proof is complete.

References

- [1] A. MAITRA and T. PARTHASARATHY, *On stochastic games*, Journ. Opti. Theory and its Appli., 5 (1970), 289-300.

- [2] K. TANAKA, S. IWASE and K. WAKUTA, *On Markov games with the expected average reward criterion*, Sci. Rep. Niigata Univ., Ser. A, No. 13 (1976), 31-41.
- [3] D. J. WHITE, *Dynamic Programming, Markov chains, and the method of successive approximations*, Journ. Math. Anal. and Appli., 6 (1963), 373-376.