

AVERAGE REWARD MARKOVIAN DECISION PROCESSES IN THE COMPLETELY ERGODIC CASE

Iwamoto, Seiichi
Department of Mathematics, Kyushu University

<https://doi.org/10.5109/13072>

出版情報：統計数理研究. 15 (3/4), pp.55-68, 1973-03. Research Association of Statistical Sciences

バージョン：

権利関係：



AVERAGE REWARD MARKOVIAN DECISION PROCESSES IN THE COMPLETELY ERGODIC CASE

By

Seiichi IWAMOTO*

(Received July 15, 1972)

1. Introduction.

We are concerned with average reward Markovian decision processes on the arbitrary state and compact action spaces. In general, Markovian decision processes are determined by four objects; S , A , q and r . S and A are non-empty Borel sets, q is a regular conditional probability on S given SA , and r is a Baire function on SAS . We interpret S as the set of states of some system, and A as the set of actions available to us at each stage. The set of actions is assumed to be independent of the state. When the system is in state s and we take action a , the system moves to a new state s' according to the conditional distribution $q(\cdot | s, a)$, and we receive a reward $r(s, a, s')$. The process is then repeated from the new state s' and so on. We then wish to maximize the expected average reward per unit time.

The purpose of this paper is to make studies on the existence of an optimal stationary policy, the policy improvement procedure under Doeblin condition and the relation between the average reward problem and the discounted total reward problem. The works of Howard [4], Derman [2] and Ross [6] were made on the problems, but in their works at least one of the state space or the action space was rather restrictive in contrast with the one in ours.

In Section 2, we shall give the basic notation and definitions to be used throughout this paper. In Section 3, assuming our processes keep up a kind of complete ergodicity resulted from Doeblin condition, we shall give sufficient conditions for the existence of an optimal stationary policy and show that a policy improvement procedure can be effectively used under some additional assumptions. In connection with the existence of an optimal policy, a functional equation fundamental to our average case will be exposed. In Section 4, we shall give the reduction of the average reward problem to the discounted total problem, and the method of successive approximations to solve the fundamental functional equation treated in Section 3. The complete ergodicity is not assumed in the reduction of the average problem to the discounted problem, but in successive approximations it is the consequence of the assumptions just set up.

* Department of Mathematics, Kyushu University, Fukuoka.

2. Probabilistic definitions and the optimization problem.

First we give the probabilistic notation and definitions following closely those of [1]. A Borel set X is a Borel subset of a Polish space. Polish space is a complete separable metric space. Unless otherwise noted, "measurable" means measurable with respect to the σ -field of a Borel subset of X , which is expressed by $\mathcal{B}(X)$. A probability is a probability measure on the measurable space $(X, \mathcal{B}(X))$, and the set of all probabilities on $(X, \mathcal{B}(X))$ is denoted by $P(X)$. If X and Y are non-empty Borel sets, a conditional probability on Y given X is a function $q(\cdot|\cdot)$ such that for each $x \in X$, $q(\cdot|x)$ is a probability on $(Y, \mathcal{B}(Y))$ and for each Borel subset B in $\mathcal{B}(Y)$, $q(B|\cdot)$ is a Baire function on X . $Q(Y|X)$ denotes the class of all conditional probabilities $q(\cdot|\cdot)$ defined as above. $\mathcal{B}(X|Y)$ will denote the set of all measurable function from X to Y . It will be shown that $\mathcal{B}(X|Y) \subset Q(Y|X)$. We denote the Cartesian product of X and Y by XY . Furthermore, the Cartesian product space of X_1, X_2, \dots will be denoted by X_1X_2, \dots , where each X_n is a Borel set. In this paper, the set of all bounded Baire functions on X is denoted by $M(X)$. $C_1(X)$ denotes the set of all bounded upper semi-continuous function of X . Obviously $C_1(X) \subset M(X)$. If $u, v \in M(X)$, $u \geq v$ means $u(x) \geq v(x)$ for all $x \in X$, and $\|u\| = \sup_{x \in X} |u(x)|$.

For any $u \in M(XY)$ and any $q \in Q(Y|X)$, qu denotes the element of $M(X)$ whose value at $x \in X$ is given by

$$qu(x) = \int_Y u(x, y) dq(y|x).$$

For any $p \in P(X)$, $q \in Q(Y|X)$, pq is the probability on XY such that for all $u \in M(XY)$, $pq(u) = p(qu)$. Any probability m on XY has a factorization $m = pq$; $p \in P(X)$ is unique and is just the marginal distribution of the first coordinate variable with respect to m , and $q \in Q(Y|X)$ is not unique and is a version of the conditional distribution of the second coordinate variable given the first.

The above notation may be extended in an obvious way to a finite or countable sequence of non-empty Borel sets X_1, X_2, \dots . If $q_n \in Q(X_{n+1}|X_1X_2 \dots X_n)$ for $n \geq 1$ and $p \in P(X_1)$, $pq_1q_2 \dots q_n$ is a probability on $X_1X_2 \dots X_{n+1}$, $pq_1q_2 \dots$ is a probability on the infinite product space $X_1X_2 \dots$, $q_2q_3 \in Q(X_3X_4|X_1X_2)$, for any $u \in M(X_1X_2 \dots X_{n+1})$, $n \geq 1$ and any m , $1 \leq m \leq n$, $q_m \dots q_n u \in M(X_1 \dots X_m)$, etc.

To avoid further complicating the notation, we shall use the following convention; for any function u on Y , we shall use the same symbol u to denote the function v on XY such that $v(x, y) = u(y)$ for all $y \in Y$. Thus, for example, if $q \in Q(Y|X)$ and $u \in M(Y)$, then $qu \in M(X)$, and $q \in Q(Y|X)$ will also denote the element $q' \in Q(Y|ZX)$ such that $q'(\cdot|z, \cdot) = q(\cdot|\cdot)$ for all $z \in Z$, etc.

A $p \in P(X)$ is degenerate if for some $x \in X$, $p\{x\} = 1$, and will sometimes be denoted by $\delta(x)$. A $q \in Q(Y|X)$ is degenerate if $q(\cdot|x)$ is degenerate for each $x \in X$. The degenerate q are exactly those which there is a measurable function f from X to Y such that $q(\cdot|x) = \delta(f(x))$ for each $x \in X$. That is to say the degenerate q is an element of $\mathcal{B}(X|Y)$. Any such $f \in \mathcal{B}(X|Y)$ will also denote its associated degenerate q , so that, for any $u \in M(XY)$, $qu(x) = u(x, f(x))$ for any $x \in X$.

Next we shall formulate our optimization problem in Markovian decision processes. Our *optimization problem* is defined by four tuple (S, A, q, r) , where S , the *set of states* of some system, is any non-empty Borel set, A , the *set of actions* available to us, is also any non-empty Borel set, $q \in Q(S|SA)$ is the *law of the motion* of the system—when the system is in s and action a is chosen, the system moves to next state s' according to the conditional distribution $q(\cdot|s, a)$; and $r \in M(SAS)$ is a *reward function*—when the system is in state s , choose action a and s' is a new state, we will receive a reward $r(s, a, s')$. We wish to maximize the expected average reward per unit time, when the process starts on state s .

A *policy* π is a sequence $\{\pi_1, \pi_2, \dots\}$, where each $\pi_n \in Q(A|H_n)$ and $H_n = SA \dots S$ ($2n-1$ factors) is the set of possible histories of the system when the n -th act must be chosen. The set of all policies is denoted by Π . A policy π is (non-randomized) *Markov* if each π_n is a degenerate element of $Q(A|S)$, i.e. $\pi = \{f_1, f_2, \dots\}$, where $f_i \in \mathcal{B}(S|A)$, and is (non-randomized) *stationary* if there is a $f \in \mathcal{B}(S|A)$ such that $\pi_n = f$ for all n . The stationary policy defined by f is denoted by $f^{(\infty)}$. The set of all stationary policies is denoted by Π_s . It is noted that $\Pi_s \subset \Pi$.

Any policy, together with law of the motion q of the system, defines for each initial state s a conditional probability distribution on the set $X = ASAS \dots$ of the future of the system, i.e. it defines an element of $Q(X|S)$, namely $e_\pi = \pi_1 q \pi_2 q \dots$. Denote the coordinate functions on SX by $s_1, a_1, s_2, a_2, \dots$, so our reward on the n -th day, as a function of the history of the system, is $r(s_n, a_n, s_{n+1})$, and total reward by the n -th day also, as a function of the history, is $\sum_{k=1}^n r(s_k, a_k, s_{k+1})$. The limiting average expected reward $J(\pi)$ from the policy π , as a function of the initial state s_1 , is then

$$J(\pi)(s_1) = \lim_{n \rightarrow \infty} \frac{1}{n} e_\pi \left[\sum_{k=1}^n r(s_k, a_k, s_{k+1}) \right].$$

It should be noted that the ambiguous notations are used.

For any $v \in M(S)$ and $0 \leq \beta \leq 1$, we denote by $I_n^\beta(\pi, v)$ the expected β -discounted total reward if we terminate after the n -th stage and receive a terminal reward $v(s_{n+1})$ at the terminal state. Thus,

$$I_n^\beta(\pi, v)(s_1) = e_\pi \left[\sum_{k=1}^n \beta^{k-1} r(s_k, a_k, s_{k+1}) + \beta^n v(s_{n+1}) \right].$$

Let $I_n^\beta(\pi) = I_n^\beta(\pi, 0)$ for $0 \leq \beta \leq 1$ and $I_n(\pi) = I_n^1(\pi)$. Then, it is clear that

$$J(\pi)(s_1) = \lim_{n \rightarrow \infty} \frac{1}{n} I_n(\pi)(s_1)$$

and that for all $\pi \in \Pi$

$$|J(\pi)| \leq \|r\|.$$

Furthermore, let $I^\beta(\pi) = \lim_{n \rightarrow \infty} I_n^\beta(\pi)$ for $0 \leq \beta < 1$. Note that this limit always exists.

3. Existence of the stationary policy and the policy improvement procedure.

In this section we are concerned with the sufficient conditions for the existence of an optimal stationary policy in the average reward Markovian decision processes on general state and action spaces. It will be shown that Howard's policy improvement can be generalized to the case of the general state and action spaces under some weak conditions introduced by Maitra [5].

Let $\{\mu_n; n \geq 1\}$ be a sequence in $P(S)$ and $\mu \in P(S)$. If $\mu_n(E)$ converges to $\mu(E)$ uniformly in $E \in \mathcal{B}(S)$ as $n \rightarrow \infty$, then we shall say that μ_n converges to μ or $\mu_n \Rightarrow \mu$ in symbol.

LEMMA 1. If $\mu_n \Rightarrow \mu$, then $\int_S g(s) d\mu_n(s)$ converges to $\int_S g(s) d\mu(s)$ as $n \rightarrow \infty$ for every $g \in M(S)$.

PROOF. The proof of this lemma is derived by extending from indicator functions.

Let $p_f(s'|s) = q(s'|s, f(s))$ for any $f \in \mathcal{B}(S|A)$, where $q(s'|s, a)$ is a given Markov transition law of the system. For any $f \in \mathcal{B}(S|A)$ n -step transition probabilities $p_f^{(n)} \equiv p_f^{(n)}(s, E)$ are easily calculated in the following:

$$p_f^{(0)} \equiv p_f^{(0)}(s, E) = \begin{cases} 1 & \text{for } s \in E \ (E \in \mathcal{B}(S)) \\ 0 & \text{otherwise,} \end{cases}$$

$$p_f^{(1)} \equiv p_f^{(1)}(s, E) = p_f(s, E),$$

$$p_f^{(n+1)} \equiv p_f^{(n+1)}(s, E) = \int_S p_f^{(n)}(s', E) p_f^{(1)}(s, ds').$$

These are also stochastic transition laws.

The following Corollary is an immediate consequence of Lemma 1.

COROLLARY. If for any fixed $f \in \mathcal{B}(S|A)$, $\frac{1}{n} \sum_{i=0}^{n-1} p_f^{(i)}(s, \cdot) \Rightarrow \pi(s, \cdot)$ for each $s \in S$ and if π is Markov transition law, then $\frac{1}{n} \sum_{i=0}^{n-1} \int_S p_f^{(i)}(s, ds') g(s')$ converges to $\int_S \pi(s, ds') g(s')$ for each $s \in S$ and any $g \in M(S)$.

Now we shall set up some hypotheses, called Doeblin condition, for any given Markov transition law $p(s, E)$ over general state space S . We shall say that the Markov transition law $p(s, E)$ satisfies Hypothesis (D) (or Doeblin condition) if there is a finite-valued measure φ of sets $A \in \mathcal{B}(S)$ with $\varphi(S) > 0$, an integer $\nu \geq 1$, and a positive ε such that

$$p^{(\nu)}(s, A) \leq 1 - \varepsilon \quad \text{if } \varphi(A) \leq \varepsilon.$$

Hypothesis (D) will be strengthened to the following Hypothesis (D'). We shall say that the Markov transition law $p(s, E)$ satisfies Hypothesis (D') if there is a measure φ on $\mathcal{B}(S)$ sets, with $0 < \varphi(S) < \infty$, an integer $\nu \geq 1$, a positive $\delta > 0$, and a $\mathcal{B}(S)$ set C for which

$$\varphi(C) > 0,$$

$$p^{(\nu)}(s, s') \geq \delta, \quad s \in S, \quad s' \in C.$$

Here $p_0^{(\nu)}(s, \cdot)$ is the density of the absolutely continuous components of $p^{(\nu)}(s, \cdot)$

with respect to φ .

We shall use the following lemma, due to Doob [3], in the policy improvement procedures.

LEMMA 2. (Doob [3]) *If the Markov transition law $p(s, E)$ satisfies Hypothesis (D'), then there is a stationary absolute probability distribution $p(\cdot) \in P(S)$, with $p(C_1) \geq \delta\varphi(C_1)$ when $C_1 \subset C$, for which*

$$|p^{(n)}(s, A) - p(A)| \leq [1 - \delta\varphi(C)]^{(n/\nu)-1}, \quad n = 1, 2, \dots$$

In the following we shall use Assumption (I): a Markov transition law resulted by any stationary policy satisfies Hypothesis (D'). That is, we shall say Assumption (I) is satisfied if for any $f \in \mathcal{B}(S|A)$ there exist, depending on f , a measure φ_f of $\mathcal{B}(S)$ sets, with $0 < \varphi_f(S) < \infty$, an integer $\nu = \nu(f) \geq 1$, a $\delta = \delta(f) > 0$, and a $\mathcal{B}(S)$ set $C = C(f)$ such that

$$\varphi(C) > 0 \quad \text{and} \quad p_0^{(\nu)}(s, s') \geq \delta, \quad s \in S, \quad s' \in C,$$

where $p_0^{(\nu)}(s, \cdot)$ is the density of the absolutely continuous components of $p_f^{(\nu)}(s, \cdot)$ with respect to φ_f . Throughout the paper, the following three assumptions given in [5] will remain operative: (1) A is a compact metric space, (2) r is a bounded upper semi-continuous function on SA and (3) $q \in Q(S|SA)$ is weak continuous, that is, $(s_n, a_n) \rightarrow (s, a)$ implies that $q(\cdot | s_n, a_n)$ converges weakly to $q(\cdot | s, a)$.

LEMMA 3. *Under Assumption (I), for any $f \in \mathcal{B}(S|A)$*

$$J(f^{(\infty)})(s) = \int_S p_f(ds') r(s', f(s')).$$

Hence average reward from the stationary policy $f^{(\infty)}$ is independent of the initial state s .

PROOF. For any $f \in \mathcal{B}(S|A)$, $\pi = \{f, f, \dots\} = f^{(\infty)} \in \Pi_s$, resulting stationary Markov transition law is $p_f(s, E) = q(E|s, f(s))$, and resulting reward is $r_f(s) = r(s, f(s))$. Hence, by definition

$$J(f^{(\infty)})(s) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{\nu=0}^{n-1} \int_S p_f^{(\nu)}(s, ds') r_f(s').$$

But Assumption (I) establishes the existence of a stationary absolute probability distribution $p_f(\cdot)$ on S independent of initial state. Furthermore, Lemma 2 yields that

$$p_f^{(\nu)}(s, E) \rightarrow p_f(E) \quad (\text{uniformly in } s \text{ and } E).$$

Hence, similarly,

$$\frac{1}{n} \sum_{\nu=0}^{n-1} p_f^{(\nu)}(s, E) \rightarrow p_f(E) \quad (\text{uniformly in } s \text{ and } E).$$

By Lemma 1 $\frac{1}{n} \sum_{\nu=1}^{n-1} \int_S p_f^{(\nu)}(s, ds') r_f(s')$ converges to $\int_S p_f(ds') r_f(s')$ for each $s \in S$. We conclude that

$$J(f^{(\infty)})(s) = \int_S p_f(ds') r_f(s').$$

This completes the proof.

In other words, Assumption (I) implies that the setting up the criteria in terms of a probability p on S , e.g. p -optimal, $p(\cdot)$ -average optimal or (p, ϵ) -optimal criteria,

is of no use. We remark that even if in Lemma 3 $r(s, a) \in C_1(SA)$ is extended to $r(s, a) \in M(SA)$, the same result remains valid.

Now we shall use the following lemmas and Selector Theorem, due to Maitra [5].

LEMMA 4 (Maitra [5]). *Let u be a bounded upper semi-continuous on SA . Define $u^*: S \rightarrow R$ by $u^*(s) = \max_{a \in A} u(s, a)$. Then u^* is upper semi-continuous on S .*

LEMMA 5 (Maitra [5]). *Let u be a bounded upper semi-continuous function on SA . Define $U: S \rightarrow 2^A$ by $U(s) = \{a \in A; u(s, a) = \max_{a' \in A} u(s, a')\}$, where 2^A is the set of all non-empty closed subset in A . Then U is a Borel map.*

SELECTOR THEOREM (Maitra [5]). *Let u be a bounded upper semi-continuous function on SA . Then there exists a Borel measurable map f from S into A such that $u(s, f(s)) = \max_{a \in A} u(s, a)$ for all $s \in S$.*

LEMMA 6 (Maitra [5]). *Let $w: S \rightarrow R$ be a bounded upper semi-continuous function. Then $g: SA \rightarrow R$ defined by $g(s, a) = \int w(\cdot) dq(\cdot | s, a)$ is upper semi-continuous on SA .*

We can now state and prove our theorem with the aid of above lemmas and Selector Theorem.

THEOREM 1. *If there exists a $\{g, v(\cdot)\}$ such that $v(\cdot) \in C_1(S)$, satisfying a functional equation*

$$g + v(s) = \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) v(s') \right\} \quad s \in S, \quad (1)$$

then there exists an $f^ \in \mathcal{B}(S|A)$ such that for $s \in S$ and every $\pi \in \Pi$*

$$g = J(f^{*(\infty)})(s) \geq J(\pi)(s).$$

Here, $f^{(\infty)}$ is the stationary policy which for each state $s \in S$ prescribes the action that maximize the right hand side of (1).*

PROOF. The proof of this theorem is similar to that of Derman [2], except for the use of above lemmas and Selector Theorem. Let a_s^* ($s \in S$) be the action which maximize the right hand side of (1). By Lemma 5 and Lemma 6, $U(s) = \{a_s^*\}$ is a Borel map from S to 2^A . Selector Theorem enables us to choose a Borel measurable $f^* \in \mathcal{B}(S|A)$ from S to A such that

$$f^*(s) \in U(s) \quad \text{for every } s \in S.$$

This reduces (1) to

$$g + v(s) = r_{f^*}(s) + \int_S p_{f^*}(s, ds') v(s'), \quad s \in S, \quad (2)$$

where

$$r_{f^*}(s) = r(s, f^*(s)) \quad \text{and} \quad p_{f^*}(s, E) = q(E | s, f^*(s)).$$

On integrating (2) by $p_{f^*}^{(v)}$, v -step transition probability, which is calculated from $p_{f^*}(s, E)$, we have

$$g + \int_S p_{f^*}^{(v)}(s_1, ds) v(s) = \int_S p_{f^*}^{(v)}(s_1, ds) r_{f^*}(s) + \int_S p_{f^*}^{(v+1)}(s_1, ds') v(s'). \quad (3)$$

An interchange of the order of the integrations in (3) is justified by Fubini's lemma.

By taking average over $\nu = 0, 1, \dots, T-1$ in (3)

$$g = \frac{1}{T} \sum_{\nu=0}^{T-1} \int_S p_{f^*}^{(\nu)}(s_1, ds) r_{f^*}(s) + \frac{1}{T} \left[v(s_1) - \int_S p_{f^*}^{(T)}(s_1, ds') v(s') \right],$$

and canceling in the limit, we have

$$\begin{aligned} g &= \lim_{n \rightarrow \infty} \frac{1}{T} \sum_{\nu=0}^{T-1} \int_S p_{f^*}^{(\nu)}(s_1, ds) r_{f^*}(s) \\ &= J(f^{*(\infty)})(s_1). \end{aligned} \quad (4)$$

Thus, g is the expected average reward per unit time from the stationary policy $f^{*(\infty)}$.

Now we must show that $f^{*(\infty)}$ is optimal over Π . We shall successively define $g_n(s)$; $n = 0, 1, 2, \dots$ by

$$g_0(s) = \max_{a \in A} r(s, a) \quad s \in S, \quad (5)$$

$$g_{n+1}(s) = \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) g_n(s') \right\} \quad s \in S.$$

Thus $g_0(\cdot)$ is well defined because of compact action space and it is clear that $g_0(\cdot) \in C_1(S)$ by Lemma 4. Lemma 4 together with Lemma 6, implies that $g_n(\cdot)$, $n = 1, 2, \dots$ are well defined and that $g_n(\cdot) \in C_1(S)$ for each n . By definition, $g_n(s)$ expresses the total expected reward obtained over the periods $0, 1, \dots, n$ operating optimally. We shall show that there exists a finite M such that for $n = 0, 1, 2, \dots$,

$$ng + v(s) - M \leq g_n(s) \leq ng + v(s) + M, \quad s \in S. \quad (6)$$

For $n = 0$, (6) holds since $v(\cdot)$ and $r(\cdot, \cdot)$ are bounded functions. Assume (6) holds for $n \leq N$. Then, by (5), (6) and (1), we get

$$\begin{aligned} g_{N+1}(s) &= \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) g_N(s') \right\} \\ &\leq \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) [Ng + v(s') + M] \right\} \\ &= Ng + M + \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) v(s') \right\} \\ &= Ng + M + g + v(s) \\ &= (N+1)g + v(s) + M, \end{aligned}$$

the right hand inequality of (6) for $N+1$. The left follows in the same way.

From the definition it follows that for any $\pi \in \Pi$, $I_n(\pi)(s)$ is the total expected reward over the periods $0, 1, \dots, n$ from π , started from state s . Since $g_n(s)$ is the results from optimal policy for those periods, we have

$$I_n(\pi)(s) \leq g_n(s), \quad n = 1, 2, \dots.$$

Hence,

$$\frac{I_n(\pi)(s)}{n} \leq \frac{g_n(s)}{n}, \quad n = 1, 2, \dots$$

Then (6) yields that

$$\liminf_{n \rightarrow \infty} \frac{I_n(\pi)(s)}{n} \leq \lim_{n \rightarrow \infty} \frac{g_n(s)}{n} = g, \quad s \in S.$$

This states that

$$J(\pi)(s) \leq g = J(f^{*(\infty)})(s), \quad s \in S,$$

which completes the proof.

COROLLARY. Under the same conditions as in Theorem 1, it holds that

$$\|g_n(s) - ng\| \leq 2M \quad \text{for } n = 1, 2, \dots$$

Next we shall consider a generalization of the policy improvement routine originally achieved by Howard [4]. Here, we set up the Assumption (II): for any $f \in \mathcal{B}(S|A)$, there exists a $\{g, v(s)\}$ such that $v(\cdot) \in C_1(S)$, satisfying the functional equation

$$g + v(s) = r(s, f(s)) + \int_S q(ds' | s, f(s))v(s'), \quad s \in S \quad (7)$$

We now define our improvement procedure under Assumptions (I) and (II). For any given $f \in \mathcal{B}(S|A)$ we construct $f' \in \mathcal{B}(S|A)$ as follows: set $f'(s) = f(s)$ for each state $s \in S$ such that

$$r(s, f(s)) + \int_S q(ds' | s, f(s))v(s') = \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a)v(s') \right\}, \quad (8)$$

and put $f'(s) = a'$ for the state s not satisfying (8), by using an action a' such that

$$r(s, a') + \int_S q(ds' | s, a')v(s') > r(s, f(s)) + \int_S q(ds' | s, f(s))v(s'). \quad (9)$$

That is,

$$f'(s) = \begin{cases} f(s) & \text{on } \{s; u(s, f(s)) = \max_{a \in A} u(s, a)\} \\ a' & \text{on } \tilde{S}(f)^c, \end{cases} = \tilde{S}(f) \in \mathcal{B}(S)$$

where

$$u(s, a) = r(s, a) + \int_S q(ds' | s, a)v(s') \in C_1(SA).$$

We remark that for example, Selector Theorem enables us to choose f' as Borel measurable, since there exists a measurable f' such that

$$u(s, f'(s)) = \max_{a \in A} u(s, a) \quad \text{on } \tilde{S}^c \in \mathcal{B}(S).$$

The following theorem shows that the policy can be properly improved in the case of general state.

THEOREM 2. Assume that Assumptions (I) and (II) hold. If $\tilde{S} = S$, then $f^{(\infty)}$ is optimal over Π . If $\tilde{S} \neq S$, then

$$J(f'^{(\infty)})(s) \geq J(f^{(\infty)})(s) \quad s \in S.$$

Furthermore if $\varphi_{f'}(\tilde{S}^c(f) \cap C(f')) > 0$, then

$$J(f'^{(\infty)})(s) > J(f^{(\infty)})(s) \quad s \in S.$$

PROOF. Let $q(s' | s, f'(s)) = p_{f'}(s, ds')$. Let $\varepsilon(s)$ be the difference between the left hand side and right one in (9); that is,

$$\varepsilon(s) = \begin{cases} r(s, f(s)) + \int_S q(ds' | s, f(s))v(s') \\ \quad - \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a)v(s') \right\} = 0 & \text{on } \tilde{S}, \\ r(s, a') + \int_S q(ds' | s, a')v(s') - r(s, f(s')) \\ \quad - \int_S q(ds' | s, f(s))v(s') > 0 & \text{on } \tilde{S}^c. \end{cases}$$

For any $s_1 \in S$ and integer ν we have, using Assumption (II), that

$$\begin{aligned} \int_S p_{f'}^{(\nu)}(s_1, ds) \varepsilon(s) &= \int_S p_{f'}^{(\nu)}(s_1, ds) r(s, f'(s)) \\ &\quad + \int_S p_{f'}^{(\nu+1)}(s_1, ds') v(s') \\ &\quad - g - \int_S p_{f'}^{(\nu)}(s_1, ds) v(s). \end{aligned} \quad (10)$$

On averaging over $\nu = 0, 1, \dots, n-1$ in (10) we get

$$\begin{aligned} \int_S \varepsilon(s) \frac{1}{n} \sum_{\nu=0}^{n-1} p_{f'}^{(\nu)}(s_1, ds) &= \frac{1}{n} \sum_{\nu=0}^{n-1} \int_S p_{f'}^{(\nu)}(s_1, ds) r_{f'}(s) \\ &\quad + \frac{1}{n} \left(\int_S p_{f'}^{(n)}(s_1, ds') v(s') - v(s_1) \right) - g. \end{aligned} \quad (11)$$

By Assumption (I), there exists the limit

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{\nu=0}^{n-1} p_{f'}^{(\nu)}(s_1, ds) = p_{f'}(ds).$$

Then on letting $n \rightarrow \infty$ in (11) we have, because of boundedness of $\varepsilon(s)$,

$$\int_S \varepsilon(s) p_{f'}(ds) = J(f'^{(\infty)})(s_1) - g.$$

On the other hand, it follows that

$$\int_S \varepsilon(s) p_{f'}(ds) = \int_{\tilde{S}^c} \varepsilon(s) p_{f'}(ds) \geq 0.$$

In particular $\varphi_{f'}(\tilde{S}^c \cap C(f')) > 0$ implies $p_{f'}(\tilde{S}^c \cap C(f')) > 0$ by Lemma 2. Thus, if $\varphi_{f'}(\tilde{S}^c \cap C(f')) > 0$, then

$$\int_{\tilde{S}^c} \varepsilon(s) p_{f'}(ds) \geq \int_{\tilde{S}^c \cap C(f')} \varepsilon(s) p_{f'}(ds) > 0,$$

namely

$$J(f^{(\infty)})(s_1) > g = J(f^{(\infty)})(s_1), \quad s_1 \in S.$$

This completes the proof of Theorem 2.

4. Reduction of average problem to discounted problem and the method of successive approximations.

We shall treat the Markovian decision processes satisfying Assumption (III): there are a state s_0 and $\alpha > 0$ such that

$$q(\{s_0\} | s, a) \geq \alpha \quad \text{for all } s \in S, a \in A.$$

For convenience we call the average reward Markovian decision process which has been stated as above the original M.D.P., and we shall define a modified M.P.D. with state space, action space and reward function not altered, but with transition probability q' altered.

In the modified M.D.P., q' is to be given by

$$q'(B | s, a) = \begin{cases} \frac{q(B | s, a)}{1 - \alpha} & \text{for } s_0 \notin B \\ \frac{q(B | s, a) - \alpha}{1 - \alpha} & \text{for } s_0 \in B \end{cases}$$

for any $B \in \mathcal{B}(S)$.

Note that q' is a Markov transition law.

We now consider the modified M.D.P. with β -discounted total expected reward over the infinite future and study its connection to the original M.D.P. with average criterion satisfying Assumption (III).

One of Maitra's results, which we shall use, is stated in

LEMMA 7 (Maitra [5]). *For any Markovian decision process (S, A, q, r) and $0 \leq \beta < 1$, there exists a stationary policy which is β -optimal over Π , and the optimal reward $(= \sup_{\pi \in \Pi} I^\beta(\pi)(s))$ is upper semi-continuous on S .*

COROLLARY. *For our modified M.D.P. (S, A, q', r) and $0 \leq \beta < 1$, there exists a β -optimal stationary policy over Π .*

It should be noted that the (s, a) -weak continuity of $q(\cdot | s, a)$ implies that of $q'(\cdot | s, a)$, since

$$\int_S g(s') dq'(s' | s, a) = \frac{1}{1 - \alpha} \int_S g(s') dq(s' | s, a) - \frac{\alpha}{1 - \alpha} g(s_0) \quad \text{for } g(\cdot) \in M(S).$$

Hence we may assume that for the modified M.D.P. (S, A, q', r) A is compact, r is bounded upper semi-continuous, and $q' \in Q(S|SA)$ is weak continuous.

Let π'_β be the β -optimal policy with respect to the modified M.D.P., and let $I'^\beta(\pi'_\beta)(s)$ be the β -discounted total expected reward with respect to the modified M.D.P., when the process starts in state s and policy π'_β is used. Note that any policy for the original M.D.P. can also be considered as a policy for the modified one and vice versa. The fundamental theorem on the reduction is the following:

THEOREM 3. Under Assumption (III), there exists an optimal stationary policy for the original M. D. P., and the optimal average reward per unit time is $\alpha I^{1-\alpha}(\pi'_{1-\alpha})(s_0)$. Furthermore, the optimal stationary policy is the one which takes the actions which maximize the right hand side of

$$\alpha I^{1-\alpha}(\pi'_{1-\alpha})(s_0) + f'(s) = \max_{a \in A} \left[r(s, a) + \int_S f'(s') dq(s' | s, a) \right], \quad (11)$$

where

$$f'(s) = I^{1-\alpha}(\pi'_{1-\alpha})(s) - I^{1-\alpha}(\pi'_{1-\alpha})(s_0). \quad (12)$$

PROOF. Blackwell [1] has shown that since $\pi'_{1-\alpha}$ is $(1-\alpha)$ -optimal with respect to the modified M. D. P. (S, A, q', r) , $I^{1-\alpha}(\pi'_{1-\alpha})(s)$ satisfies the optimal equation, namely,

$$I^{1-\alpha}(\pi'_{1-\alpha})(s) = \max_{a \in A} \left[r(s, a) + (1-\alpha) \int_S I^{1-\alpha}(\pi'_{1-\alpha})(s') dq'(s' | s, a) \right], \quad s \in S.$$

By (12) we have

$$\begin{aligned} f'(s) + I^{1-\alpha}(\pi'_{1-\alpha})(s_0) &= (1-\alpha) I^{1-\alpha}(\pi'_{1-\alpha})(s_0) \\ &\quad + \max_{a \in A} \left[r(s, a) + (1-\alpha) \int_S f'(s') dq'(s' | s, a) \right]. \end{aligned}$$

This yields that

$$\alpha I^{1-\alpha}(\pi'_{1-\alpha})(s_0) + f'(s) = \max_{a \in A} \left[r(s, a) + \int_S f'(s') dq(s' | s, a) \right]. \quad (13)$$

Thus, $\{\alpha I^{1-\alpha}(\pi'_{1-\alpha})(s_0), f'(s)\}$ ($f'(\cdot) \in C_1(S)$) satisfies the optimal equation with respect to the average case, that is, condition of Theorem 1 is satisfied. It follows that there exists an optimal stationary policy $f^{*(\infty)} \in \Pi_s$ with respect to the original M. D. P. i. e.

$$\alpha I^{1-\alpha}(\pi'_{1-\alpha})(s_0) = J(f^{*(\infty)})(s) \geq J(\pi)(s) \quad s \in S, \quad \pi \in \Pi,$$

and that $f^{*(\infty)}$ is the stationary one which takes the actions that maximize the right hand side of (11). This completes the proof of the Theorem 3.

We now show that the solution $\{g, v(\cdot)\}$ of the functional equation (1) can be solved by the method of successive approximations under Assumption (III)': there exist an integer $\nu \geq 0$, a quantity $0 < \alpha \leq 1$, and a state $s_0 \in S$ such that for all $a_1, a_2, \dots, a_{\nu+1} \in A$ and $s \in S$

$$q^{(\nu+1)}(\{s_0\} | s; a_1, a_2, \dots, a_{\nu+1}) \geq \alpha$$

holds. This is a generalization of White's method. It should be noted that Assumption (III) implies Assumption (III)'.

THEOREM 4. Under Assumption (III)', the sequence $\{g_n, v_n(\cdot), n \geq 1\}$ defined by

$$\begin{aligned} V_n(s) &= \max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) v_{n-1}(s') \right\}, \quad s \in S, \\ g_n &= V_n(s_0), \\ v_n(s) &= V_n(s) - g_n, \end{aligned} \quad (14)$$

converges uniformly to the solution $\{g, v(\cdot)\}$ of the functional equation (1),

$$g+v(s)=\max_{a \in A} \left\{ r(s, a) + \int_S q(ds' | s, a) v(s') \right\} \quad s \in S, \quad (1)$$

where

$$v_0(s) = \max_{a \in A} r(s, a).$$

PROOF. By Lemma 4 and Lemma 6, it follows that $V_n(\cdot)$, $v_n(\cdot)$ are well-defined and that $V_n(\cdot)$, $v_n(\cdot) \in C_1(S)$. We need only prove the uniformity of convergence, since this implies that the limiting form is a solution of (1). For any sequence of upper semi-continuous functions $\{V_n(s)\}$ define

$$\begin{aligned} \nabla_n(V) &= \inf_{s \in S} [V_n(s) - V_{n-1}(s)] \\ \Delta_n(V) &= \sup_{s \in S} [V_n(s) - V_{n-1}(s)], \end{aligned}$$

and similarly,

$$\begin{aligned} \nabla_n(v) &= \inf_{s \in S} [v_n(s) - v_{n-1}(s)] \\ \Delta_n(v) &= \sup_{s \in S} [v_n(s) - v_{n-1}(s)]. \end{aligned}$$

When $n \geq \nu + 3$, it can be easily shown that

$$\begin{aligned} V_n(s) - V_{n-1}(s) &\geq \inf_{a \in A} \int_S q(ds' | s, a) [v_{n-1}(s') - v_{n-2}(s')] \\ &= \inf_{a \in A} \left\{ \int_S q(ds' | s, a) [V_{n-1}(s') - V_{n-2}(s')] - (V_{n-1}(s_0) - V_{n-2}(s_0)) \right\}. \end{aligned}$$

Hence,

$$\begin{aligned} &V_n(s) - V_{n-1}(s) + \{V_{n-1}(s_0) - V_{n-2}(s_0)\} \\ &\geq \inf_{a \in A} \left\{ \int_S q(ds' | s, a) \left[\inf_{a' \in A} \int_S q(ds'' | s', a') [V_{n-2}(s'') - V_{n-3}(s'')] \right] \right. \\ &\quad \left. - \{V_{n-2}(s_0) - V_{n-3}(s_0)\} \right\}, \end{aligned}$$

namely,

$$\begin{aligned} &V_n(s) - V_{n-1}(s) + \{V_{n-1}(s_0) - V_{n-2}(s_0)\} + \{V_{n-2}(s_0) - V_{n-3}(s_0)\} \\ &\geq \inf_{a \in A} \left\{ \int_S q(ds' | s, a) \left[\inf_{a' \in A} \int_S q(ds'' | s', a') [V_{n-2}(s'') - V_{n-3}(s'')] \right] \right\} \\ &= \inf_{a, a'} \int_S q^{(2)}(ds'' | s; a, a') [V_{n-2}(s'') - V_{n-3}(s'')]. \end{aligned}$$

Repeating this iteration procedure $(\nu - 1)$ times, we have

$$\begin{aligned} &V_n(s) - V_{n-1}(s) + \sum_{k=1}^{\nu} \{V_{n-k}(s_0) - V_{n-k-1}(s_0)\} \\ &\geq \inf_{a_1, a_2, \dots, a_{\nu}} \int_S q^{(\nu)}(s' | s; a_1, a_2, \dots, a_{\nu}) \{V_{n-\nu}(s') - V_{n-\nu-1}(s')\}, \quad \text{i. e.} \\ &V_n(s) - V_{n-1}(s) + \{V_{n-1}(s_0) - V_{n-\nu-1}(s_0)\} \\ &\geq \inf_{a_1, a_2, \dots, a_{\nu+1}} \int_S q^{(\nu+1)}(s' | s; a_1, a_2, \dots, a_{\nu+1}) \{v_{n-\nu-1}(s') - v_{n-\nu-2}(s')\}. \end{aligned}$$

Now by hypothesis

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

and

$$q(\{s_0\} | s; a_1, a_2, \dots, a_{\nu+1}) \geq \alpha,$$

then we have

$$\mathcal{V}_n(V) + \{V_{n-1}(s_0) - V_{n-\nu-1}(s_0)\} \geq (1-\alpha)\mathcal{V}_{n-\nu-1}(v). \quad (15)$$

Similarly,

$$\mathcal{A}_n(V) + \{V_{n-1}(s_0) - V_{n-\nu-1}(s_0)\} \leq (1-\alpha)\mathcal{A}_{n-\nu-1}(v). \quad (16)$$

(15) and (16) imply that

$$\mathcal{A}_n(V) - \mathcal{V}_n(V) \leq (1-\alpha)(\mathcal{A}_{n-\nu-1}(v) - \mathcal{V}_{n-\nu-1}(v)). \quad (17)$$

By the way,

$$\mathcal{V}_n(v) = \mathcal{V}_n(V) - (g_n - g_{n-1})$$

and

$$\mathcal{A}_n(v) = \mathcal{A}_n(V) - (g_n - g_{n-1})$$

yield

$$\mathcal{A}_n(v) - \mathcal{V}_n(v) = \mathcal{A}_n(V) - \mathcal{V}_n(V). \quad (18)$$

From (17), (18)

$$\mathcal{A}_n(v) - \mathcal{V}_n(v) \leq (1-\alpha)\{\mathcal{A}_{n-\nu-1}(v) - \mathcal{V}_{n-\nu-1}(v)\}.$$

Let $D_n(v) = \mathcal{A}_n(v) - \mathcal{V}_n(v)$. Then we have for some finite C , and for $n = N(\nu+1) + r$ ($0 \leq r \leq \nu$),

$$D_n(v) \leq (1-\alpha)^N D_r(v) \leq (1-\alpha)^N C.$$

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

$$\mathcal{V}_n(v) \leq 0 \leq \mathcal{A}_n(v).$$

If

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

and

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

then

$$U_n(v) \leq D_n(v) \leq (1-\alpha)^N C.$$

Hence $\{v_n(s), n \geq 1\}$ converges uniformly to the function $v(s)$.

On the other hand, $U_n(V) \leq U_{n-1}(v)$ implies the uniform convergence of $\{V_n(s)\}$ to the function $V(s)$, so does the sequence $\{v_n(s_0)\} = \{g_n\}$. The interchange of the order of maximum and limit can be justified by the inequalities

$$|\max_{a \in A} f_n(s, a) - \max_{a \in A} f(s, a)| \leq \sup_{a \in A} |f_n(s, a) - f(s, a)|$$

and

$$\|f_n - f\| \leq \|v_{n-1} - v\|,$$

where

$$f_n(s, a) = r(s, a) + \int_S q(s' | s, a) v_{n-1}(s') \in C_1(SA)$$

and

$$f(s, a) = r(s, a) + \int_S q(s' | s, a) v(s') \in C_1(SA).$$

Acknowledgement: The author would like to express his deep appreciation to Doctor N. Furukawa whose advices and inspiration helped make his paper possible.

References

- [1] D. BLACKWELL, *Discounted dynamic programming*, Ann. Math. Statist., **36** (1965), 226-235.
- [2] C. Derman, *Denumerable state Markovian decision processes-average cost criterion*, Ann. Math. Statist., **37** (1966), 1545-1554.
- [3] J. L. DOOB, *Stochastic Processes*, Wiley, New York (1953).
- [4] R. A. HOWARD, *Dynamic Programming and Markov Processes*, Wiley, New York (1960).
- [5] A. MAITRA, *Discounted dynamic programming on compact metric spaces*, Sankhyā Ser A., **30** (1968), 211-216.
- [6] S. M. ROSS, *Arbitrary state Markovian Decision Processes*, Ann. Math. Statist., **39** (1968), 2118-2122.
- [7] D. J. WHITE, *Dynamic Programming, Markov chains, and the method of successive approximations*, J. Math. Anal. Appl., **6** (1963), 373-376.