

## MARKOV DECISION PROCESSES WITH A NEW OPTIMALITY CRITERION: DISCRETE TIME<sup>1</sup>

BY STRATTON C. JAQUETTE

*Cornell University*

Standard finite state and action discrete time Markov decision processes with discounting are studied using a new optimality criterion called moment optimality. A policy is moment optimal if it lexicographically maximizes the sequence of signed moments of total discounted return with a positive (negative) sign if the moment is odd (even). This criterion is equivalent to being a little risk adverse. It is shown that a stationary policy is moment optimal by examining the negative of the Laplace transform of the total return random variable. An algorithm to construct all stationary moment optimal policies is developed. The algorithm is shown to be finite.

**1. Introduction.** This paper is concerned with finite state and action discrete time Markov decision processes where future returns are discounted. We define and study a new optimality criterion which we call moment optimality.

The Markov decision process as well as most of the notation needed to develop the results in this paper are defined in Section 2. We also define and discuss moment optimality. Section 3 contains the result that there is a stationary policy that is optimal under the criterion of moment optimality. As the methods of Section 3 do not yield a workable algorithm to construct optimal policies, a slightly different approach is given in Section 4 to obtain an algorithm to construct all stationary moment optimal policies. Section 5 is devoted to showing that the algorithm is finite.

**2. Preliminaries.** We consider a standard Markov decision process. We assume that the stochastic process has a finite state space denoted by  $S$ , and without loss of generality assume that  $S = \{1, 2, \dots, s\}$ . We assume that the process starts at time  $t = 0$  and that the process can jump from state to state at discrete points in time  $t = 0, 1, 2, \dots$ . At each point in time an action is selected and applied to the Markov process. We assume that there is a single finite set of actions, which we denote by  $A$ . Any element  $a$  in  $A$  is an action which may be applied to the process in any state and at any point in time. We can just as easily assume that there are perhaps distinct finite action sets  $A_i$  available when the process is in state  $i$  and obtain the same results, but we make the more restrictive assumption to simplify the exposition.

If the process is in state  $i$  at some point in time, applying action  $a$  determines the return and transition probabilities for that period. The return obtained is

---

Received May 1971; revised September 1972.

<sup>1</sup> The original work for this paper was partially supported by Office of Naval Research Contract N00014-67-A-0112-0031 (NR-042-265) and National Science Foundation Grant GP-8790.

denoted by  $r(i, a)$ . The returns are bounded since both state space and action sets are finite. We denote the transition probabilities by  $p_{ij}(a)$  ( $j = 1, 2, \dots, s$ ). The transition probability  $p_{ij}(a)$  is the probability that the process will be in state  $j$  at time  $t + 1$  given that the process is in state  $i$  and action  $a$  is applied at time  $t$ . We assume that  $p_{ij}(a) \geq 0$  for all  $i, j$ , and  $a$ , and that  $\sum_{j=1}^s p_{ij}(a) = 1$  for all  $i$  and  $a$ .

We define the set  $F$  by  $F = \prod_{i=1}^s A$ . The elements of  $F$ , denoted  $\mathbf{f}$ , are called action vectors. The  $i$ th component of  $\mathbf{f}$ , denoted  $f(i)$ , is the action taken if the process is in state  $i$ .

We denote a policy by  $\pi$  and can write this policy as a sequence of action vectors, e.g.  $\pi = \{\mathbf{f}_0, \mathbf{f}_1, \dots, \mathbf{f}_t, \dots\}$ . Here  $\mathbf{f}_t$  is the action vector applied at time  $t$  using policy  $\pi$ , and we can write  $\pi(t) = \mathbf{f}_t$ . If  $\pi$  is stationary we have  $\pi(0) = \pi(t)$  ( $t = 1, 2, \dots$ ), and if  $\pi(0) = \mathbf{f}$ , we usually write  $f^\infty$  to denote this stationary policy.

The state of the process is a random variable which depends on the starting state of the process, the policy used, and the time in question. We denote by  $X_\pi(t)$  the random state of the system at time  $t$  when policy  $\pi$  is used and suppress the dependence of  $X_\pi(t)$  on the starting state.

A return is obtained in each period. This return is a function of the random variable  $X_\pi(t)$  and is itself therefore a random variable. In explicit notation the return at time  $t$  when policy  $\pi$  is used is  $r(X_\pi(t), \mathbf{f}(X_\pi(t)))$  where  $\mathbf{f} = \pi(t)$ . We shall simplify this to  $r(X_\pi(t), \pi)$ . Future returns are discounted by a constant discount factor  $\beta$  ( $0 \leq \beta < 1$ ) per period.

We denote by  $\mathbf{r}(f)$  the column vector of returns associated with the action vector  $\mathbf{f}$ . We denote by  $P(f)$  the one period transition probability matrix associated with  $\mathbf{f}$ . The vector  $\mathbf{R}(\pi)$  will denote the total discounted return using policy  $\pi$ . The total discounted return random variable given that  $X_\pi(0) = i$  is:

$$(2.1) \quad \begin{aligned} [\mathbf{R}(\pi)]_i &= r(i, \pi) + \beta r(X_\pi(1), \pi) + \beta^2 r(X_\pi(2), \pi) + \dots \\ &= \sum_{t=0}^{\infty} \beta^t r(X_\pi(t), \pi) \end{aligned} \quad (X_\pi(0) = i).$$

Let  $\mathbf{X}_\pi(t)$  be the vector whose  $i$ th component is 1 if  $X_\pi(t) = i$  and is 0 otherwise and let  $\circ$  denote componentwise multiplication i.e.,  $[\mathbf{u} \circ \mathbf{v}]_i \equiv [\mathbf{u}]_i [\mathbf{v}]_i$  if  $\mathbf{u}$  and  $\mathbf{v}$  have the same dimension. Then we can rewrite (2.1) as

$$(2.2) \quad \mathbf{R}(\pi) = \sum_{t=0}^{\infty} \beta^t \mathbf{r}(\pi(t)) \circ \mathbf{X}_\pi(t).$$

We can characterize the total discounted return random vector  $\mathbf{R}(\pi)$  by its moments. We denote by  $\mathbf{M}_n(\pi)$  the  $n$ th moment of  $\mathbf{R}(\pi)$ . Noting that  $\mathbf{R}(\pi)^n = \mathbf{R}(\pi) \circ \mathbf{R}(\pi) \circ \dots \circ \mathbf{R}(\pi)$ , we define  $\mathbf{M}_n(\pi)$  as follows:

$$(2.3) \quad \begin{aligned} \mathbf{M}_n(\pi) &\equiv E[\mathbf{R}(\pi)^n] & (n = 1, 2, \dots) \\ \mathbf{M}_0(\pi) &\equiv \mathbf{1}. \end{aligned}$$

We also consider  $\mathbf{M}_n(\pi)$  multiplied by  $+1$  or  $-1$  depending on whether  $n$  is odd

or even. For this reason we define the vectors  $\mathbf{N}_n(\pi)$  as follows:

$$(2.4) \quad \mathbf{N}_n(\pi) \equiv (-1)^{n+1} \mathbf{M}_n(\pi) \quad (n = 0, 1, 2, \dots).$$

We also use vectors  $\mathbf{N}^m(\pi)$  whose  $m + 1$  components are themselves vectors. We define  $\mathbf{N}^m(\pi)$  as follows:

$$(2.5) \quad \mathbf{N}^m(\pi) \equiv (\mathbf{N}_0(\pi), \mathbf{N}_1(\pi), \dots, \mathbf{N}_m(\pi)) \quad (m = 0, 1, 2, \dots).$$

By  $\mathbf{N}(\pi)$  we mean  $\mathbf{N}^\infty(\pi)$ . We use the negative of the Laplace transform of the total discounted return random vector extensively and denote this by  $\mathbf{U}_\pi(\lambda)$ :

$$(2.6) \quad \mathbf{U}_\pi(\lambda) \equiv -E[\exp(-\lambda \mathbf{R}(\pi))].$$

By expression (2.6) we mean that  $[\mathbf{U}_\pi(\lambda)]_i = -E[\exp(-\lambda [\mathbf{R}(\pi)]_i)]$ .

Suppose  $\mathbf{u}$  and  $\mathbf{v}$  are two vectors of the same dimension with components  $u(i)$  and  $v(i)$  respectively. We write  $\mathbf{u} > \mathbf{v}$ , that  $\mathbf{u}$  is lexicographically greater than  $\mathbf{v}$ , if there is an integer  $n$  such that  $u(i) = v(i)$  for  $i < n$  and  $u(n) > v(n)$ . If  $\mathbf{N}$  and  $\mathbf{M}$  are vectors whose components are themselves vectors, then we write  $\mathbf{N} > \mathbf{M}$  to signify that there is an integer  $n$  such that  $\mathbf{N}(i) = \mathbf{M}(i)$  for  $i < n$  and  $\mathbf{N}(n) > \mathbf{M}(n)$ . The relationships  $<$ ,  $\leq$ , and  $\geq$  are defined analogously. In this paper we shall write  $\mathbf{u} < \mathbf{v}$  for vectors if  $\mathbf{u} \leq \mathbf{v}$  and  $\mathbf{u} \neq \mathbf{v}$ .

The optimality criterion which has been used almost exclusively is maximization of the expected value of total return. Exceptions to this are Derman [1] and Whitt [4], both for the case of infinite horizon and no discounting: Derman considers the sample path average return and Whitt considers weak convergence of the total return suitably normalized. We choose instead to retain the fact that the total return obtained from using a given policy is a random variable and pose the question of what sort of random return would be desirable. Using the notation above we say that a policy  $\pi^*$  is *moment optimal* if  $\mathbf{N}(\pi^*) \geq \mathbf{N}(\pi)$  for all policies  $\pi$ .

Moment optimality measures the desirability of a policy by examining the moments of the associated return random variable. A policy is good if its odd moments are large and even moments small, where we consider the moments lexicographically. Thus if two policies have unequal means, the one with the greater mean is better, as in the usual expected value case. If they have equal means, the one with the smaller variance is better. If they have equal means and variances, the one with the greater third moment is better, and so forth. In some sense a return distribution is good if it has large expected value with small risk of falling far below the expected value.

Moment optimality will not differ from the usual expected value case if there is a unique stationary policy which maximizes the expected return. In general there may be several stationary policies which attain the maximum expected return and many more nonstationary policies which do as well. Examples are easy to construct; a simple two state example can be found in Jaquette ([3], page 15).

Moment optimality can also be viewed as being a "little bit risk adverse"

when an exponential utility function is used. If the utility function  $-e^{-\lambda x}$  is assumed for the utility of having wealth  $x$ , then  $\lambda$  is the constant aversion to risk implicit in the utility function. The problem of finding a policy which yields maximum expected utility of total discounted return for this utility function when the aversion to risk is small enough is equivalent to the problem of finding a moment optimal policy. This will be clear in the treatment which follows, although it will not be stated explicitly. For a discussion of this and other implications for utility theory and Markov decision processes, the reader should see Jaquette ([3] Chapter 5), and Howard and Matheson ([2]). The approach here and in Jaquette ([3]) is different from that used by Howard and Matheson; however the results are quite similar.

**3. A stationary policy is moment optimal.** The main result of this section is:

**THEOREM 1.** *For the finite state and action discrete time Markov decision process defined above, there exists a stationary policy which is moment optimal.*

We dispense with some preliminaries before proving this theorem. We know that  $r(\cdot, \cdot)$  is bounded; assume that  $r(i, a) \leq B < \infty$  for all  $i \in S$  and  $a \in A$ . From (2.1) we conclude that  $|\mathbf{R}(\pi)|_i \leq (1 - \beta)^{-1}B < \infty$ . Thus we can conclude that  $\mathbf{R}(\pi)$  is bounded everywhere and that its moments grow at most geometrically since  $|\mathbf{M}_n(\pi)| \leq B^n \mathbf{1}$ . This is sufficient to ensure that the Laplace transform of  $\mathbf{R}(\pi)$ , and hence  $\mathbf{U}_\pi(\lambda)$ , exists and is finite on the whole real line for any policy  $\pi$ . We are thus justified in using the Taylor series expansion for  $\mathbf{U}_\pi(\lambda)$ .

From the definition of the function  $\mathbf{U}$  we have

$$[\mathbf{U}_{f\pi}(\lambda)]_i = -E[\exp(-\lambda[\mathbf{R}(f\pi)]_i)],$$

where by the policy  $f\pi$  we mean that  $\mathbf{f}$  is used at time zero and then the elements of  $\pi$  are used in sequence. The Markov property of the process allows us to conclude that

$$(3.1) \quad [\mathbf{U}_{f\pi}(\lambda)]_i = \exp(-\lambda[\mathbf{r}(f)]_i) \sum_{j=1}^s P_{ij}(f(i))[\mathbf{U}_\pi(\beta\lambda)]_j.$$

This result follows directly by using the Markov property and taking suitable conditional expectations. Equation (3.1) can be written in vector and matrix form as follows:

$$(3.2) \quad \mathbf{U}_{f\pi}(\lambda) = \{\exp(-\lambda\mathbf{r}(f))\} \circ \{P(f)\mathbf{U}_\pi(\beta\lambda)\}.$$

We can simplify the calculations by defining the operator  $L_f$  and restating (3.2) in terms of this operator. Let  $\mathcal{L}$  be the space of functions of a real variable which are the negative of Laplace transforms of bounded random variables. Let  $\mathcal{L}^s$  be the  $s$ -fold direct product of  $\mathcal{L}$ :  $\mathcal{L}^s = \prod_{i=1}^s \mathcal{L}$ . We define the operator  $L_f$  as follows:

$$(3.3) \quad (L_f \mathbf{u})(\lambda) \equiv \{\exp(-\lambda\mathbf{r}(f))\} \circ \{P(f)\mathbf{u}(\beta\lambda)\}, \quad \mathbf{u}(\lambda) \in \mathcal{L}^s.$$

It is easy to verify that  $L_f \mathbf{u}$  is in  $\mathcal{L}^s$  if  $\mathbf{u}$  is in  $\mathcal{L}^s$  and thus that  $L_f: \mathcal{L}^s \rightarrow \mathcal{L}^s$ . We can then restate (3.2) as:

$$(3.4) \quad \mathbf{U}_{f\pi}(\lambda) = L_f \mathbf{U}_\pi(\lambda).$$

The following lemma establishes a sense in which  $L_f$  is monotone.

LEMMA 3.1. *Let  $\mathbf{f}$  be any action vector,  $\lambda_0 > 0$ , and  $\mathbf{u}$  in  $\mathcal{L}^s$ . If*

$$(3.5) \quad L_f(\mathbf{u})(\lambda) \geq \mathbf{u}(\lambda)$$

for all  $\lambda$  in  $[0, \lambda_0]$ , then for all  $\lambda$  in  $[0, \lambda_0]$

$$\mathbf{U}_{f^\infty}(\lambda) \geq \mathbf{u}(\lambda).$$

PROOF. Equations (3.3) and (3.5) suffice to establish inductively that for  $\lambda$  in  $[0, \lambda_0]$

$$(L_f)^n(\mathbf{u})(\lambda) \geq \mathbf{u}(\lambda).$$

It suffices to show that  $(L_f)^n(\mathbf{u})(\lambda) \rightarrow \mathbf{U}_{f^\infty}(\lambda)$ . To show this define a policy  $\pi^n \pi^*$  to be the policy that uses the action vectors from  $\pi$  until time  $n$  and then uses the action vectors from  $\pi^*$  in their proper order. If  $\pi = \{\mathbf{f}_0, \mathbf{f}_1, \dots\}$ , define  $L_\pi^n$  as follows:

$$L_\pi^n(\mathbf{u}) \equiv L_\pi^{n-1}(L_{f_n}(\mathbf{u})), \quad \text{and} \quad L^0(\mathbf{u}) \equiv L_{f_0}(\mathbf{u}).$$

It is elementary to show that  $\mathbf{R}(\pi^n \pi^*) \rightarrow \mathbf{R}(\pi)$  a.e. ( $n \rightarrow \infty$ ), that  $\mathbf{U}_{\pi^n \pi^*}(\lambda) \rightarrow \mathbf{U}_\pi(\lambda)$  ( $n \rightarrow \infty$ ), and that

$$(3.6) \quad L_\pi^n(\mathbf{u}) \rightarrow \mathbf{U}_\pi \quad (n \rightarrow \infty) \quad \text{on the real line.}$$

We note that Lemma 3.1 and its proof hold with all inequalities reversed. It also holds with a strict inequality. We shall use these generalizations in what follows.

LEMMA 3.2. *Let  $\mathbf{f}$  and  $\mathbf{g}$  be any action vectors in  $F$ . Define the function  $h(\cdot)$  as follows:*

$$(3.7) \quad h(\lambda) \equiv [\mathbf{U}_{f^\infty}(\lambda)]_i - [\mathbf{U}_{g^\infty}(\lambda)]_i.$$

There exists a positive number  $\lambda_0$  such that  $h(\cdot)$  does not change sign on the interval  $[0, \lambda_0]$ .

PROOF. Laplace transforms are analytic on the interior of their region of convergence. Thus  $h(\cdot)$  is analytic on the real line. If the lemma were false, then  $h$  would cross zero infinitely often as  $\lambda$  approaches zero, and there would exist a sequence  $\{\lambda_n\} \rightarrow 0^+$  such that  $h(\lambda_n) = 0$  for all  $n$ . This would imply that  $h$  is identically zero, since  $h$  is analytic. Thus the lemma must hold.

For convenience in notation for the following lemma, we write  $h > 0$  ( $\lambda \rightarrow 0^+$ ) if a function as defined in Lemma 3.2 does not change sign on some interval  $[0, \lambda_0]$  and if  $h(\lambda) > 0$  on  $(0, \lambda_0)$ .

LEMMA 3.3. Let  $\mathbf{f}$  and  $\mathbf{g}$  be any action vectors in  $F$  and let  $h$  be as in Lemma 3.2. Define the action vector  $\mathbf{e}$  as follows:

$$(3.8) \quad \begin{aligned} e(i) &\equiv f(i) && \text{if } h > 0 \quad (\lambda \rightarrow 0^+), \\ &\equiv g(i) && \text{otherwise.} \end{aligned}$$

There exists a scalar  $\lambda_0 > 0$  such that

$$U_{\mathbf{e}^\infty}(\lambda) \geq \max \{U_{\mathbf{f}^\infty}(\lambda), U_{\mathbf{g}^\infty}(\lambda)\} \quad \text{for all } \lambda \in [0, \lambda_0].$$

PROOF. Associated with each index  $i$  there is a positive  $\lambda_0^i$  guaranteed by Lemma 3.2 such that the appropriate  $h(\cdot)$  does not change sign. Take the  $\lambda_0$  for Lemma 3.3 to be the minimum of these  $\lambda_0^i$ , which must be positive since  $S$  is finite. The definition for  $\mathbf{e}$ , (3.8), leads directly to the statements that for  $\lambda \in [0, \lambda_0]$

$$L_{\mathbf{e}}(U_{\mathbf{f}^\infty})(\lambda) \geq U_{\mathbf{f}^\infty}(\lambda), \quad \text{and} \quad L_{\mathbf{e}}(U_{\mathbf{g}^\infty})(\lambda) \geq U_{\mathbf{g}^\infty}(\lambda).$$

We can now apply Lemma 3.1 to obtain the conclusion of the lemma.

We now define an alternative optimality criterion. A policy  $\pi^*$  is called *U-optimal* if there exists a positive scalar  $\lambda_0$  such that  $U_{\pi^*}(\lambda) \geq U_\pi(\lambda)$  for every  $\pi$  and all  $\lambda$  in  $[0, \lambda_0]$ . We remark that *U-optimality* is optimality being a “little bit risk adverse.”

LEMMA 3.4. There exists a stationary policy that is *U-optimal*.

PROOF. Restrict consideration first to stationary policies. Starting with any stationary policy and applying Lemma 3.3 repeatedly, we are assured of the existence of a positive  $\lambda_0$  and an action vector  $\mathbf{f}$  such that  $U_{\mathbf{f}^\infty}(\lambda) \geq U_{\mathbf{g}^\infty}(\lambda)$  for all  $\lambda \in [0, \lambda_0]$  and for every  $\mathbf{g} \in F$ . This follows since there are only a finite number of distinct elements in  $F$ . We can then conclude immediately that every action vector  $\mathbf{g}$  satisfies

$$(3.9) \quad L_{\mathbf{g}}(U_{\mathbf{f}^\infty})(\lambda) \leq U_{\mathbf{f}^\infty}(\lambda) \quad \text{for all } \lambda \in [0, \lambda_0].$$

If this conclusion did not follow, then Lemma 3.3 would ensure an improvement on  $\mathbf{f}$ , which cannot obtain. By suitable choices of  $\mathbf{g}$  in (3.9) and applying these  $L_{\mathbf{g}}$  repeatedly, we can construct  $L_\pi^n$  for any choice of policies  $\pi$ . This application ensures that for all  $\lambda \in [0, \lambda_0]$

$$L_\pi^n(U_{\mathbf{f}^\infty})(\lambda) \leq U_{\mathbf{f}^\infty}(\lambda).$$

Equation (3.6) allows us to conclude that  $U_\pi(\lambda) \leq U_{\mathbf{f}^\infty}(\lambda)$  on  $[0, \lambda_0]$ , which completes the proof.

LEMMA 3.6. A policy is *moment optimal* if and only if it is *U-optimal*.

PROOF. This follows directly from the definitions of moment optimality and *U-optimality*, Lemma 3.5, and the Taylor series expansion for  $U_\pi(\lambda)$ , which converges for all  $\lambda$ ,

$$(3.10) \quad U_\pi(\lambda) = \sum_{n=0}^{\infty} N_n(\pi) \lambda^n / n!.$$

It is clear the  $N(\pi^*) \geq N(\pi)$  if and only if  $U_{\pi^*}(\lambda) \geq U_{\pi}(\lambda)$  for all  $\lambda$  is some suitable interval  $[0, \lambda_0]$  ( $\lambda_0 > 0$ ).

Theorem 1 can now be verified very easily. Lemma 3.5 establishes the existence of a stationary policy that is  $U$ -optimal. Lemma 3.6 establishes that the moment optimal policies are exactly the  $U$ -optimal policies.

It should also be clear from the proof that the assumption of a single action set  $A$  is merely for convenience and that the more general action sets  $A_i$  with  $F = \prod_{i=1}^s A_i$  admits the identical proofs.

It should also be clear that the particular selection of alternating signs in the definition of moment optimality is not critical in the proofs. In fact the same results obtain if we choose any of the following sequence of signs for the moments of return in the definition of  $N(\pi)$  in (2.4):  $+1, -1, +1, \dots$  or  $-1, +1, -1, \dots$  or  $+1, +1, +1, \dots$  or  $-1, -1, -1, \dots$ .

**4. An algorithm to construct a moment optimal policy.** To facilitate construction of all stationary policies which are moment optimal, we define additional optimality criteria which are related to moment optimality. We define the criterion of ( $m$ ) *moment optimality* to be moment optimality ignoring all moments higher than the  $m$ th moment. Recalling the definition of  $N^m(\pi)$  in (2.5), we call a policy  $\pi^*$  ( $m$ ) *moment optimal* if  $N^m(\pi^*) \geq N^m(\pi)$  for all policies  $\pi$ .

It should be evident that a policy is moment optimal if and only if it is ( $m$ ) moment optimal for all  $m$ . We base an algorithm on this observation. From Theorem 1 we know that there is a stationary policy which is moment optimal. A stationary moment optimal policy is ( $m$ ) moment optimal for all  $m$ . We can, therefore, restrict our attention to stationary policies.

Define the sets of action vectors,  $F^m$ , as follows:

$$(4.1) \quad F^m \equiv \{f: f \in F \text{ and } f^\infty \text{ is } (m) \text{ moment optimal}\}.$$

From our previous observations  $F^m \supseteq F^{m+1}$  and  $F^\infty = \lim_{m \rightarrow \infty} F^m \neq \emptyset$ . The set  $F^\infty$  is the set of action vectors,  $f$ , such that  $f^\infty$  is moment optimal.

In this section we shall use some additional notation. If  $N$  and  $M$  are vectors with vector components  $N_i$  and  $M_i$  respectively, we write  $N =_n M$  if  $N_i = M_i$  for  $i \leq n$ . Similarly we write  $N \geq_n M$  if  $N > M$  and  $N \neq_n M$ . The relation  $\geq_n, \leq_n$ , and  $<_n$  are similarly defined.

Suppose that  $U$  is in  $\mathcal{L}^s$ ; then we can write the Taylor expansion as  $U(\lambda) = \sum_{n=0}^\infty u_n \lambda^n/n!$ . By the notation  $U$  we mean the vector whose vector components are given by this Taylor expansion, e.g.,  $U = \{u_0, u_1, \dots\}$ . If  $U(\lambda)$  in  $\mathcal{L}^s$  is viewed as the negative of the Laplace transform of some bounded random vector  $R$ , then  $E[R^n] = (-1)^{n+1}u_n$ .

LEMMA 4.1. *Suppose  $U(\lambda)$  is in  $\mathcal{L}^s$  and let  $f$  be any action vector. Then*

$$(4.2) \quad U_{f^\infty} =_m U$$

*if and only if*

$$(4.3) \quad L_f(U) =_m U.$$

Suppose that  $L_f(\mathbf{U}) \succ_m \mathbf{U}$ . Then

$$\mathbf{U}_{f^\infty} \succ_m \mathbf{U}.$$

PROOF. Suppose that (4.3) holds. Then a simple induction argument leads to  $(L_f)^n(\mathbf{U}) =_m \mathbf{U}$ . The result in (3.6) concerning convergence of the transform implies  $(L_f)^n(\mathbf{U}) \rightarrow \mathbf{U}_{f^\infty}$ , which yields (4.2).

Suppose that (4.2) holds. Using (3.4) establishes that  $L_f(\mathbf{U}_{f^\infty}) = \mathbf{U}_{f^\infty}$ , therefore if we can show that  $L_f(\mathbf{U}) =_m L_f(\mathbf{U}_{f^\infty})$  follows from (4.2) we could establish that

$$L_f(\mathbf{U}) =_m L_f(\mathbf{U}_{f^\infty}) = \mathbf{U}_{f^\infty} =_m \mathbf{U},$$

which is just (4.3). To fill the gap expand the components of (3.3) as a power series in  $\lambda$  and equate coefficients of  $\lambda^n$ :

$$(4.4) \quad [L_f(\mathbf{U})]_n = \sum_{i=0}^n \beta^i \binom{n}{i} (-\mathbf{r}(f))^{n-i} \circ P(f) \mathbf{u}_i.$$

Since  $[L_f(\mathbf{U})]_n$  is a linear combination of  $\mathbf{u}_i$  for  $i \leq n$  only and since  $\mathbf{U}_{f^\infty} =_m \mathbf{U}$  by (4.2), we can conclude that  $[L_f(\mathbf{U})]_n = [L_f(\mathbf{U}_{f^\infty})]_n$  for  $n \leq m$ . This completes this portion of the proof.

If  $L_f(\mathbf{U}) \succ_m \mathbf{U}$ , then using (3.10) it is clear that the hypotheses of Lemma 3.1 are satisfied and hence that  $\mathbf{U}_{f^\infty} \geq \mathbf{U}$ . Applying  $L_f$  to both sides of this relationship yields

$$\mathbf{U}_{f^\infty} = L_f(\mathbf{U}_{f^\infty}) \geq L_f(\mathbf{U}) \succ_m \mathbf{U}.$$

LEMMA 4.2. *The set  $F^m$  is a direct product. If  $f \in F^m$  and  $F^m = \times_{i=1}^s A_i^m$ , then*

$$(4.5) \quad A_i^m = \{a \in A_i : L_g(\mathbf{U}_{f^\infty}) =_m \mathbf{U}_{f^\infty} \text{ for } g(i) = a \text{ and } g(j) = f(j) \ (i \neq j)\}.$$

PROOF. Let  $\mathbf{U}$  in Lemma 4.1 be  $\mathbf{U}_{f^\infty}$ . Since  $\mathbf{f} \in F^m$ , Lemma 4.1 implies that  $\mathbf{g} \in F^m$  if and only if  $L_g(\mathbf{U}_{f^\infty}) =_m \mathbf{U}_{f^\infty}$ . Clearly  $F^{m-1} \supseteq F^m$ , so that to define  $F^m$  we need only consider  $\mathbf{g} \in F^{m-1}$ . This eliminates consideration of all components of  $L_g(\mathbf{U}_{f^\infty})$  except the  $m$ th. From (4.4) it is evident that the  $i$ th component of  $[L_g(\mathbf{U})]_m$  depends only on  $g(i)$  and not on  $g(j)$  for  $j \neq i$ . It thus follows directly that  $F^m$  is a direct product and that  $A_i^m$  can be given by (4.5).

We remark that  $A_i^{m-1} \supseteq A_i^m$  follows directly from  $F^{m-1} \supseteq F^m$ . Since  $F^\infty$  is nonempty by Theorem 1 so are the  $A_i^m$ . This indicates that  $F^m$  can easily be constructed from  $F^{m-1}$  if an element  $\mathbf{f}$  in  $F^m$  can be found. This construction can be accomplished by using (4.5) and (4.4) and noting that  $[\mathbf{U}_{f^\infty}]_m = \mathbf{N}_m(f^\infty)$ .

LEMMA 4.3. *Pick any  $\mathbf{f}$  in  $F^{m-1}$ . Either*

- (a)  $\mathbf{f} \in F^m$  or
- (b) *there exists a  $\mathbf{g} \in F^{m-1}$  such that*

$$(4.6) \quad L_g(\mathbf{U}_{f^\infty}) \succ_m \mathbf{U}_{g^\infty}.$$

If (4.6) holds, then  $\mathbf{U}_{g^\infty} \succ_m \mathbf{U}_{f^\infty}$ .

PROOF. First assume that (4.6) holds. Lemma 4.1 can be applied to show directly that  $\mathbf{U}_{g^\infty} \succ_m \mathbf{U}_{f^\infty}$ .



Suppose then that (4.6) holds for no  $\mathbf{g}$  in  $F^{m-1}$ . For every  $\mathbf{g}$  in  $F^{m-1}$  we then have  $L_g(\mathbf{U}_{f^\infty}) \leq_m \mathbf{U}_{f^\infty}$ . A direct analog of Lemma 4.1 for the inequality  $<$  and that lemma as stated combine to show that then  $\mathbf{U}_{g^\infty} \leq_m \mathbf{U}_{f^\infty}$  for all  $\mathbf{g}$  in  $F^{m-1}$ . This is simply the statement that  $\mathbf{f}$  is in  $F^m$ , which is part (a) of the lemma.

Lemma 4.3 indicates a policy improvement algorithm. An action vector  $\mathbf{g}$  is sought which improves upon all others in  $F^{m-1}$  in the sense that it is in  $F^m$ . Lemma 4.1, especially (4.4) indicates how to find a  $\mathbf{g}$  in  $F^{m-1}$  satisfying (4.6). Since  $L_g(\mathbf{U}_{f^\infty}) =_{m-1} \mathbf{U}_{f^\infty}$ , we seek an improvement  $\mathbf{g}$  satisfying  $[L_g(\mathbf{U}_{f^\infty})]_m > [\mathbf{U}_{f^\infty}]_m$  or

$$(4.7) \quad \sum_{i=0}^m \beta^i \binom{m}{i} (-\mathbf{r}(g))^{m-i} \circ P(g) \mathbf{N}_i(f^\infty) > \mathbf{N}_m(f^\infty).$$

The  $i$ th component of the left side of (4.7) depends only on  $g(i)$  and not on  $g(j)$  ( $j \neq i$ ). Thus an improvement on  $\mathbf{f}$  can be obtained by  $\mathbf{g}$  defined by

$$g(i) = \{a \in A_i^{m-1} : a \text{ maximizes the } i\text{th component of the left side of (4.7)}\}.$$

If  $\mathbf{g}$  obtained equals  $\mathbf{f}$ , then  $\mathbf{f}$  is in  $F^m$  and no further improvement is possible. This  $\mathbf{f}$  then allows direct construction of the  $A_i^m$  and  $F^m$  as indicated by Lemma 4.2 and the following remarks.

**THEOREM 2.** *The set  $F^m$  can be constructed in a finite number of steps.*

**PROOF.** The theorem follows by an inductive argument using Lemmas 4.2 and 4.3 and the following remarks. The finiteness of each step, generating  $F^m$  from  $F^{m-1}$ , follows since the policy improvement can be performed only a finite number of times. This is a result of the finiteness of  $F$ .

**5. The algorithm to construct moment optimal policies is finite.** Theorem 2 assures that ( $m$ ) moment optimal policies can be constructed a finite algorithm. Moment optimality requires consideration of all moments, and construction of ( $\infty$ ) moment optimal policies would require a countable algorithm. In fact only a finite algorithm is needed, as indicated in the following theorem.

**THEOREM 3.** *Moment optimality is equivalent to ( $n^*$ ) moment optimality for some finite  $n^*$ . All stationary moment optimal policies can be constructed by a finite algorithm.*

**PROOF.** Since there are only a finite number of stationary policies, there must be a finite number  $n^*$  such that for all pairs of action vectors  $\mathbf{f}$  and  $\mathbf{g}$ ,  $\mathbf{U}_{f^\infty} > \mathbf{U}_{g^\infty}$  if and only if  $\mathbf{U}_{f^\infty} >_{n^*} \mathbf{U}_{g^\infty}$ .

Under certain circumstances  $n^*$  can be determined beforehand to be no greater than  $s$ ; however in general this is not the case. If the algorithm is performed iteratively, a stopping rule can be given: if  $\mathbf{R}(f^\infty) =_{\mathscr{D}} \mathbf{R}(g^\infty)$  for all  $\mathbf{f}$  and  $\mathbf{g}$  in  $F^m$  for some  $m$ , then  $n^* \leq m$ . This is true if  $F^m$  has a single element; however there may be several elements in  $F^\infty$  with the same return distributions, in which case verifying that  $\mathbf{R}(f^\infty) =_{\mathscr{D}} \mathbf{R}(g^\infty)$  is equivalent to calculating a countable number of moments.

**6. Acknowledgments.** The author would like to acknowledge the assistance of his advisor, Professor Donald L. Iglehart, throughout the preparation of his dissertation while at Stanford University. These results are given in much greater detail in the dissertation [3]. The author also appreciates the suggestions of the referee, which were very helpful in indicating a more compact method of proof.

## REFERENCES

- [1] DERMAN, C. (1964). On sequential control processes. *Ann. Math. Statist.* **35** 341-349.
- [2] HOWARD, R. A. and MATHESON, J. E. (1972). Risk-sensitive Markov decision processes. *Management Sci.* **18** 356-369.
- [3] JAQUETTE, S. C. (1971). Markov decision processes with a new optimality criterion. Technical Report No. 15, Department of Operations Research, Stanford Univ.
- [4] WHITT, W. (1972). Stochastic abelian and tauberian theorems. *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete* **22** 251-267.

CORNELL UNIVERSITY  
DEPARTMENT OF OPERATIONS RESEARCH  
UPSON HALL  
ITHACA, NEW YORK 14850