

## CONDITIONS FOR THE EQUIVALENCE OF OPTIMALITY CRITERIA IN DYNAMIC PROGRAMMING

BY JAMES FLYNN

University of Chicago

This paper examines the relationships between optimality criteria which are commonly used for undiscounted, discrete-time, countable state Markovian decision models. One approach, due to Blackwell, is to maximize the expected discounted total return as the discount factor approaches 1. Another, due to Veinott, is to maximize the Cesaro means of the finite horizon expected returns as the horizon tends to infinity. Derman's is to maximize the long-run average gain. Denardo, Miller and Lippman showed that Blackwell's and Veinott's approaches are equivalent for finite state and action spaces. As shown here, that equivalence breaks down when the state space is countable. Also, policies optimal according to Blackwell's or Veinott's approach need not be optimal according to Derman's. On the positive side, fairly weak conditions are given under which Blackwell's and Veinott's criteria imply Derman's, and somewhat stronger conditions under which Blackwell's and Veinott's criteria are equivalent.

**1. Introduction.** Our formulation of the Markovian decision model follows Blackwell's (1965). Consider a system with a countable *state space*  $S$  and a finite *action space*  $A$ . Each day the current state  $s \in S$  is observed and an action  $a \in A$  is selected. This results in (1) an *immediate income*  $i(s, a)$  and (2) a transition to a new state  $s'$  with probability  $q(s' | s, a)$ . The incomes are assumed bounded. The problem is to control the system in the most effective manner over an infinite future. A rule or *policy*  $\pi$  for controlling the system specifies for each  $n \geq 1$  what act to choose on the  $n$ th day as a function of the system's current history  $h = (s_1, a_1, \dots, s_n)$  or, more generally,  $\pi$  specifies for each  $h$  a probability distribution on  $A$ . A (nonrandomized) *stationary policy* is a policy which is specified by a single function  $f$  mapping  $S$  into  $A$ : under it, you select act  $f(s)$  whenever the system is in state  $s$ . Given an initial state  $s$  and a policy  $\pi$ , let  $r_j(s, \pi)$  denote the *expected return* on the  $j$ th day ( $j = 1, 2, \dots$ ); then the *expected  $N$ -stage return* is  $V^N(s, \pi) = \sum_{j=1}^N r_j(s, \pi)$ , the *average gain* in the first  $N$  periods is  $V^N(s, \pi)/N$ , and the *expected discounted total return* is  $V_\beta(s, \pi) = \sum_{j=1}^\infty \beta^{j-1} r_j(s, \pi)$  where  $0 \leq \beta < 1$  is the discount factor.

For the discounted problem, there is only one reasonable criterion: a policy  $\pi_*$  is said to be  $\beta$ -optimal if

$$(1) \quad V_\beta(s, \pi_*) = V_\beta^*(s) \equiv \sup_\pi V_\beta(s, \pi), \quad s \in S.$$

Received June 1974; revised February 1976.

AMS 1970 subject classifications. Primary 49C15, 62L99, 90C40, 93C55; Secondary 60J10, 60J20.

Key words and phrases. Dynamic programming, Markovian decision process, optimality criteria, average overtaking criteria, average gain, discounting, small interest rates.

Blackwell (1965) and Maitra (1965) established the existence of a  $\beta$ -optimal stationary policy for each  $0 \leq \beta < 1$ .

There are, however, a number of competing criteria for the undiscounted problem. Derman (1964), (1966), and Ross (1968) used criteria which depend on the average gain:

DEFINITION 1. A policy  $\pi_*$  is *lim sup average optimal* if

$$(2) \quad \limsup_n V^n(s, \pi_*)/n \geq \limsup_n V^n(s, \pi)/n, \quad s \in S$$

for every policy  $\pi$ ; it is *lim inf average optimal* if

$$(3) \quad \liminf_n V^n(s, \pi_*)/n \geq \liminf_n V^n(s, \pi)/n, \quad s \in S$$

for every  $\pi$ ; and it is *average optimal* if

$$(4) \quad \liminf_n (V^n(s, \pi_*) - V^n(s, \pi))/n \geq 0, \quad s \in S$$

for every  $\pi$ . Given any  $\varepsilon > 0$ , the policy  $\pi_*$  is  $\varepsilon$ -*average optimal* if

$$(5) \quad \liminf_n (V^n(s, \pi_*) - V^n(s, \pi))/n \geq -\varepsilon, \quad s \in S$$

for every  $\pi$ .

Unfortunately, those criteria are underselective since they depend only on the tail of the returns and not on the returns during the first millenium. One can, of course, get around this with the criterion of selecting a  $\pi_*$  such that  $\liminf_N (V^N(s, \pi_*) - V^N(s, \pi)) \geq 0$  for all  $s \in S$  and all  $\pi$ . However, the latter is overselective since even when  $S$  is finite, there need not exist any policy satisfying it (see Denardo and Miller (1968)). One can circumvent both problems by using the following criterion, which is due to Veinott (1966):

DEFINITION 2. A policy  $\pi_*$  is *average overtaking optimal* if

$$(6) \quad \liminf_N \sum_{n=1}^N (V^n(s, \pi_*) - V^n(s, \pi))/N \geq 0, \quad s \in S$$

for every policy  $\pi$ .

Another way of approaching the undiscounted problem is to study the discounted problem for the case of *small interest rates* (i.e., values of  $\beta$  close to 1). The following criteria are due to Blackwell (1962) (see Remark 2 at the end of this section).

DEFINITION 3. A policy  $\pi_*$  is *optimal* if there exists a  $\beta_0 \in (0, 1)$  such that

$$(7) \quad V_\beta(s, \pi_*) \geq V_\beta(s, \pi), \quad s \in S, \beta \in (\beta_0, 1)$$

for every policy  $\pi$ .

DEFINITION 4. A policy  $\pi_*$  is *1-optimal* if

$$(8) \quad \liminf_{\beta \rightarrow 1^-} (V_\beta(s, \pi_*) - V_\beta(s, \pi)) \geq 0, \quad s \in S$$

for every policy  $\pi$ .

Optimality certainly implies 1-optimality. Using an Abelian argument (see

Hobson (1926)), one can establish the following general result:

$$(9) \quad \begin{array}{l} \text{optimality} \Rightarrow \\ \text{average overtaking optimality} \Rightarrow \end{array} 1\text{-optimality}.$$

It is also clear that

$$(10) \quad \text{average optimality} \begin{array}{l} \Leftarrow \text{lim sup average optimality} \\ \Leftarrow \text{lim inf average optimality} \end{array}.$$

One is tempted to link (9) and (10) by the statement: “1-optimality implies average optimality.” Surprisingly, this is not the case. In Flynn (1974), we give an example where an optimal policy is not lim inf average optimal. (We also show that this “pathology” cannot occur when  $S$  is finite.) Here we give examples of policies which are both optimal and average overtaking optimal but are not lim inf average optimal and lim sup average optimal, respectively (see Section 5, Examples 3 and 4). We also investigate:

QUESTION 1. When does

$$(11) \quad 1\text{-optimality} \Rightarrow \text{average optimality?}$$

Denardo and Miller (1968) and Lippman (1969) established the following:

**THEOREM 1.** *If  $S$  is finite, then average overtaking optimality and 1-optimality are equivalent.*

This paper contains an example which shows that optimality need not imply average overtaking optimality (see Section 5, Example 1). Hence, Theorem 1 does not extend to countable  $S$ . Also, Blackwell (1962) gave an example (with  $S$  finite) which indicates that neither 1-optimality nor average overtaking optimality imply optimality. It follows that (9) cannot be strengthened without imposing some conditions.

QUESTION 2. When is it true that

$$(12) \quad \text{average overtaking optimality} \Leftrightarrow 1\text{-optimality?}$$

The following related—but different (see Ross (1971))—questions are of interest:

QUESTION 1'. When does 1-optimality among stationary policies imply average optimality among stationary policies?

QUESTION 2'. When is average overtaking optimality among stationary policies equivalent to 1-optimality among stationary policies?

To give this second group of questions proper meaning, we require:

**DEFINITION 5.** Let  $\Omega$  be a class of policies and let  $\varepsilon > 0$  be a real number. A policy  $\pi_*$  is *average overtaking optimal*, *optimal*, *1-optimal*,  *$\varepsilon$ -average optimal*, or *average optimal among  $\pi$  in  $\Omega$*  if (6), (7), (8), (5) or (4), respectively, holds for

all  $\pi \in \Omega$ . When  $\Omega$  denotes the class of stationary policies, replace the phrase "among  $\pi$  in  $\Omega$ " by "among stationary policies."

REMARK 1. One can easily show (9) holds if we replace each type of optimality by the corresponding type of optimality "among  $\pi \in \Omega$ " where  $\Omega$  is arbitrary. Also, Theorem 1 remains valid if we replace each type by the corresponding type "among stationary policies."

REMARK 2. Our definition of "1-optimality" (see Definition 4) is equivalent to Veinott's (1969) definition of "0-discount optimality" but differs slightly from his definition of "1-optimality." Veinott (1966) calls  $\pi_*$  "1-optimal" if it satisfies

$$(13) \quad \lim_{\beta \rightarrow 1^-} (V_\beta(s, \pi_*) - V_\beta^*(s)) = 0, \quad s \in S.$$

(This is the same as Blackwell's (1962) definition of "nearly optimal.") Our definitions are equivalent when  $S$  is finite since there always exists an optimal policy in that case (see Blackwell (1962)). However, for general  $S$ , average overtaking optimality would not imply 1-optimality if 1-optimality were defined via (13) (see Section 5, Example 7).

**2. Outline and discussion.** Our results on Questions 1 and 1' appear in Section 3, while those on Questions 2 and 2' appear in Section 4. All of our examples are in Section 5. The proofs of some of our technical results are relegated to Section 6. In this section, we focus on some of our more striking results.

Our answer to Question 1 is fairly complete. In Corollary 2 we show that the following condition is sufficient for (11): For every  $\epsilon > 0$ , there exists a policy which is  $\epsilon$ -average optimal. (This result is fairly useful: see Remark 3 below.) Surprisingly, the analogous result does not hold for Question 2. In Example 2, we construct a model in which there exists an average overtaking optimal policy and a 1-optimal policy which is not average overtaking optimal. Apparently, the relationship between 1-optimality and average overtaking optimality is not all that strong. It would be interesting to answer:

QUESTION 3. Is the existence of an average overtaking optimal policy sufficient for optimality to imply average overtaking optimality?

REMARK 3. By Corollary 2, the existence of an average optimal policy suffices for 1-optimality to imply average optimality. Now, many authors (e.g., Derman (1964), (1966), Derman and Veinott (1967), and Ross (1968), (1971)) have investigated the question of the existence of a stationary lim inf average optimal policy; for the most part, their results apply directly to the problem of the existence of an average optimal policy (see Theorem 2 below). That condition (iv) of Theorem 2 would be enough to eliminate the "pathology" described in Flynn (1974) was originally conjectured by Bennett Fox.

THEOREM 2. *There exists a stationary policy which is average optimal in each of the following cases:*

- (i)  $S$  is finite.

(ii) *There exists a bounded set of numbers  $\{g, f(s)\}$ ,  $s \in S$ , satisfying*

$$(14) \quad g + f(s) = \min_{a \in A} \{i(a, s) + \sum_{s' \in S} q(s' | s, a)f(s')\} \quad s \in S .$$

*In this case, any stationary policy  $\pi_*$  which, for each  $s$ , selects the action which minimizes the RHS of (14) is average optimal.*

(iii)  $V_{\beta}^*(s') - V_{\beta}^*(s'')$  *is bounded uniformly in  $\beta$ ,  $s'$  and  $s''$ .*

(iv) *There exists a family  $\{\pi^{\beta}, 0 \leq \beta < 1\}$  of  $\beta$ -optimal stationary policies and a state  $s_0$  with the property: under the Markov chain associated with any  $\pi^{\beta}$ , the system eventually reaches  $s_0$  with probability 1; moreover, the mean recurrence time from  $s \in S$  to  $s_0$  under  $\pi_{\beta}$  is bounded uniformly in  $\beta$  and  $s$ .*

PROOF. The sufficiency of (i) follows from Theorem 4.2 of Brown (1965); that of (ii) follows from Remark 1 of Ross (1968), who also showed that both (iii) and (iv) imply (ii).

In our investigation of Questions 1' and 2', we determine conditions under which the following are true (see Remark 4 below):

STATEMENT 1. Any stationary policy  $\pi_*$  which is 1-optimal among stationary policies is average optimal among stationary policies.

STATEMENT 2. Any stationary policy  $\pi_*$  which is 1-optimal among stationary policies is average overtaking optimal among stationary policies.

In Theorem 4, we show that the following condition implies Statement 1: (EPR) *Under the Markov chain associated with each stationary policy the system eventually reaches a positive recurrent state with probability 1.* Note that (EPR) holds when  $S$  is finite (Chung (1967)).

In Theorem 5, we find conditions which imply Statement 2. In order to give an example where these conditions apply, we introduce the notion of *uniform recurrence*. Let  $C$  be an indecomposable class of states in a Markov chain with  $n$ -step transition matrix  $P^{(n)} = (P_{ij}^{(n)})$ , ( $i, j, n = 1, 2, \dots$ ). We say that  $C$  is *uniformly recurrent* if  $C$  is closed and if for some  $j \in C$  there exists a  $\delta > 0$  and a positive integer  $n$  such that  $\sum_{k=1}^n P_{ij}^{(k)} \geq \delta$  for all  $i \in C$ . One can show that a uniformly recurrent class of states consists of a nonempty set of positive recurrent states and a (possibly empty) set of transient states. The expected time spent in the transient states is bounded. (Uniform recurrence is equivalent to the notion of uniform  $\phi$ -recurrence which appears in Orey (1971) when  $\phi$  is the measure which assigns unit mass to  $j$  and zero mass to every other state.) Theorem 6 states that the following condition implies Statement 2: (UR) *Under the Markov chain associated with each stationary policy, the state space can be expressed as a union of uniformly recurrent classes.* (Observe that the classes need not be disjoint.) Note that (UR) does not include the case where  $S$  is finite. (Statement 1 is, of course, always true in that case: see Remark 1.)

REMARK 4. If a (possibly nonstationary) policy  $\pi_*$  is optimal among stationary policies, then under (EPR) it is average optimal among stationary policies

while under (UR) it is average overtaking optimal among stationary policies (see Corollaries 3 and 4). It would be interesting to answer:

**QUESTION 4.** Under (EPR), need a (possibly nonstationary) 1-optimal policy be average optimal among stationary policies?

Note that the answer to the analogous question for average overtaking optimal policies is "no" (see Example 8).

**REMARK 5.** None of the positive results in Sections 3 and 4 require that  $A$  be finite; however, some of them fail if we relax the condition that income be bounded (see Examples 5 and 6).

**3. Questions 1 and 1'.** This section deals with the relationship between 1-optimality and average optimality. Corollary 2 summarizes our main result on Question 1, while Theorem 4 and Corollary 3 summarize our main results on Question 1'. The next result is our key.

**THEOREM 3.** *Let  $\Omega$  be any class of policies such that for each positive  $\varepsilon$ ,  $\Omega$  contains a policy which is  $\varepsilon$ -average optimal among  $\pi$  in  $\Omega$ . Then any  $\pi_* \in \Omega$  which is 1-optimal among  $\pi$  in  $\Omega$  must also be average optimal among  $\pi$  in  $\Omega$ .*

The following corollaries are immediate consequences of Theorem 3.

**COROLLARY 1.** *If  $\lim_n V^n(s, \pi)/n$  exists for all  $s \in S$  and  $\pi \in \Omega$ , then any  $\pi_* \in \Omega$  which is 1-optimal among  $\pi$  in  $\Omega$  must also be average optimal among  $\pi$  in  $\Omega$ .*

**COROLLARY 2.** *If for every  $\varepsilon > 0$  there exists a policy which is  $\varepsilon$ -average optimal, then 1-optimality implies average optimality.*

**REMARK 6.** By Corollary 1, the following convergence condition implies Statement 1: (C)  $\lim_n V^n(s, \pi)/n$  exists for all  $s \in S$  and all stationary  $\pi$ . Note that (C) does not imply the existence of stationary  $\varepsilon$ -average optimal policies (see Ross (1971)). In Lemma 2 below we show that (EPR) implies (C), hence the following:

**THEOREM 4.** *Condition (EPR) implies Statement 1.*

We shall show that Theorem 4 implies:

**COROLLARY 3.** *Under (EPR), if  $\pi_*$  is optimal among stationary policies, then it is average optimal among stationary policies.*

We do not know whether the analogue of Corollary 3 holds for 1-optimal  $\pi_*$  (see Remark 4 above). Note that both Theorem 4 and Corollary 3 fail if we relax the boundedness condition on  $i(\cdot, \cdot)$  (see Example 5).

Theorem 3 requires the following lemma, the proof of which is relegated to Section 6.

**LEMMA 1.** *For every  $M > 0$  and  $\varepsilon > 0$ , there exists a  $\delta > 0$  such that*

$$(15) \quad \liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{j=1}^{\infty} \beta^{j-1} a_j - \liminf_n \sum_{j=1}^n a_j/n < \varepsilon$$

for any sequence  $\{a_j\}_{j=1}^\infty$  satisfying

$$(16) \quad |a_j| \leq M, \quad j = 1, 2, \dots,$$

and

$$(17) \quad \limsup_n \sum_{j=1}^n a_j/n - \liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{j=1}^\infty \beta^{j-1} a_j < \delta.$$

PROOF OF THEOREM 3. We can assume that there exists a policy  $\pi_*$  which is 1-optimal among  $\pi$  in  $\Omega$ . By hypothesis, there exist policies  $\pi_1, \pi_2, \dots$  satisfying

$$(18) \quad \liminf_n \sum_{j=1}^n [r_j(s, \pi_k) - r_j(s, \pi)]/n \geq -1/k$$

for every  $\pi \in \Omega$ , ( $k = 1, 2, \dots$ ). In particular,

$$(19) \quad \limsup_n \sum_{j=1}^n [r_j(s, \pi_*) - r_j(s, \pi_k)]/n \leq 1/k, \quad k = 1, 2, \dots.$$

The fact that  $\pi_*$  is 1-optimal among  $\pi$  in  $\Omega$  implies

$$(20) \quad \liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{j=1}^\infty \beta^{j-1} (r_j(s, \pi_*) - r_j(s, \pi)) \geq 0$$

for every  $\pi \in \Omega$ . By (19) and (20), the fact that  $i(\cdot, \cdot)$  is bounded, and Lemma 1, there exists a subsequence  $\{\pi_{m_k}\}_{k=1}^\infty$  of policies satisfying

$$(21) \quad \liminf_n \sum_{j=1}^n (r_j(s, \pi_*) - r_j(s, \pi_{m_k}))/n \geq -1/k, \quad k = 1, 2, \dots.$$

Thus, (18) and (21) imply that

$$\begin{aligned} \liminf_n \sum_{j=1}^n (r_j(s, \pi_*) - r_j(s, \pi))/n &\geq \liminf_n \sum_{j=1}^n (r_j(s, \pi_*) - r_j(s, \pi_{m_k}))/n \\ &\quad + \liminf_n \sum_{j=1}^n (r_j(s, \pi_{m_k}) - r_j(s, \pi))/n \\ &\geq -1/k - 1/m_k \end{aligned}$$

for every policy  $\pi \in \Omega$ , ( $k = 1, 2, \dots$ ). Since  $k$  is arbitrary, (4) holds. The theorem follows.

Theorem 4 requires the following lemma (see Remark 6).

LEMMA 2. Condition (EPR) implies condition (C).

PROOF. Let the random variables  $X_1, X_2, \dots$  represent the Markov chain associated with a given stationary policy  $\pi$  and initial state  $s_0 \in S$ . Define the function  $f$  on  $S$  by  $f(s) = i(s, \pi(s))$ , ( $s \in S$ ). Observe that  $f$  is bounded. This and condition (EPR) imply that  $(f(X_1) + \dots + f(X_n))/n$  converges almost everywhere to a bounded random variable  $R$  (see Chung (1967): 1.15, Theorem 2). By the bounded convergence theorem,  $(Ef(X_1) + \dots + Ef(X_n))/n$  converges to  $ER$ . To finish the proof of the lemma, note that  $Ef(X_n) = r_n(s, \pi)$ , ( $n = 1, 2, \dots$ ).

PROOF OF COROLLARY 3. Let  $\pi_*$  be optimal among stationary policies. Since there exist  $\beta$ -optimal stationary policies for  $0 \leq \beta < 1$  (see Blackwell (1965)),  $\pi_*$  must be optimal. Using this together with the arguments in the proof of Theorem 8.3 of Strauch (1966), one can show there exists an optimal stationary policy  $\pi_{**}$ ; moreover,

$$(22) \quad r_n(s, \pi_*) = r_n(s, \pi_{**}), \quad s \in S; n = 1, 2, \dots.$$

Corollary 3 follows easily from (22) and Theorem 4.

**4. Question 2 and 2'.** This section deals with the relationship between 1-optimality and average overtaking optimality. Our positive results on Questions 2 and 2' are less impressive than those on Questions 1 and 1'; in particular, the analogues of Theorem 3 and Corollary 2 are false: the existence of an average overtaking optimal policy does not imply any equivalence between 1-optimality and average overtaking optimality (see Example 2). However, we do have the following analogue of Corollary 1.

**THEOREM 5.** *Let  $\Omega$  be a class of policies such that  $g(s, \pi) \equiv \lim_n V^n(s, \pi)/n$  and  $b(s, \pi) \equiv \lim \sum_{j=1}^N (V^j(s, \pi) - jg(s, \pi))/N$  exist and are finite for all  $s \in S$  and all  $\pi \in \Omega$ . If  $\pi_* \in \Omega$  and  $\pi_*$  is 1-optimal among  $\pi \in \Omega$ , then  $\pi_*$  is average overtaking optimal among  $\pi \in \Omega$ .*

**REMARK 7.** In Section 3, we showed that (C) implies Statement 1 (see Remark 6). Unfortunately, (C) does not imply Statement 2 (see Example 1). By Theorem 5, Statement 2 holds under the following convergence condition which is stronger: (CS)  $g(s, \pi) \equiv \lim_n V^n(s, \pi)/n$  and  $b(s, \pi) \equiv \lim_N \sum_{j=1}^N (V^j(s, \pi) - jg(s, \pi))/N$  exist and are finite for all  $s \in S$  and all stationary  $\pi$ . This condition holds when  $S$  is finite (see Doob (1953)). Note that (EPR) does not imply (CS) (counterexamples are easy to construct). In Lemma 3 below, we show that (UR) implies (CS), hence the following:

**THEOREM 6.** *Condition (UR) implies Statement 2.*

Arguing as in the proof of Corollary 3, one can establish:

**COROLLARY 4.** *Under (UR), if  $\pi_*$  is optimal among stationary policies, then it is average overtaking optimal among stationary policies.*

The analogue of Corollary 4 does not hold for 1-optimal  $\pi_*$ : under (UR), a nonstationary policy can be 1-optimal without being average overtaking optimal (see Example 8). Note that both Theorem 6 and Corollary 4 fail if we relax the boundedness condition on  $i(\cdot, \cdot)$  (see Example 6).

**PROOF OF THEOREM 5.** Let  $\pi_* \in \Omega$  be 1-optimal among policies in  $\Omega$  and let  $\pi$  be an arbitrary policy in  $\Omega$ . Fix  $s \in S$ . By Corollary 1,  $g(s, \pi_*) \geq g(s, \pi)$ . One can easily show that  $g(s, \pi_*) > g(s, \pi)$  implies  $\lim_N \sum_{n=1}^N (V^n(s, \pi_*) - V^n(s, \pi)) = \infty$ , which is stronger than (6). Thus we need only consider the case  $g(s, \pi_*) = g(s, \pi)$ . In this case,  $\liminf_{\beta \rightarrow 1^-} (V_\beta(s, \pi_*) - V_\beta(s, \pi)) = \liminf_{\beta \rightarrow 1^-} \{ \sum_{j=1}^\infty \beta^{j-1} (r_j(s, \pi_*) - g(s, \pi_*)) - \sum_{j=1}^\infty \beta^{j-1} (r_j(s, \pi) - g(s, \pi)) \}$ . By condition (CS) and an Abelian argument (see Hobson (1926)), the RHS equals

$$\lim_N \sum_{n=1}^N \sum_{j=1}^n (r_j(s, \pi_*) - g(s, \pi_*))/N - \lim_N \sum_{n=1}^N \sum_{j=1}^n (r_j(s, \pi) - g(s, \pi))/N.$$

Since  $g(s, \pi_*) = g(s, \pi)$ , the latter equals  $\lim_N \sum_{n=1}^N \sum_{j=1}^n (r_j(s, \pi_*) - r_j(s, \pi))/N$ . This and the 1-optimality of  $\pi_*$  give us (6).

Theorem 6 requires only the following lemma (see Remark 7).

**LEMMA 3.** *Condition (UR) implies condition (CS).*



PROOF. Let  $\pi$  be a stationary policy. Fix  $s \in S$ . Since (UR) implies (EPR), the existence of  $\lim_n V^n(s, \pi)/n$  follows from Lemma 2. Using Theorem 7.1 of Orey (1971) and the interpretation of uniform  $\phi$ -recurrence given in Section 2, one can show that  $\lim_N \sum_{j=1}^N (V^j(s, \pi) - j \lim_n V^n(s, \pi)/n)/N$  exists and is finite. We are done.

**5. Counterexamples.** This section contains all of our counterexamples. In Example 1, there exists a stationary optimal policy which is not average overtaking optimal. In the second example, there exists a stationary average overtaking optimal policy, and a stationary 1-optimal policy which is not average overtaking optimal. (Hence, the existence of an average overtaking optimal policy is not sufficient for (12).) In the third and fourth, there exist stationary policies which are both optimal and average overtaking optimal but are not  $\liminf$  average optimal and  $\limsup$  average optimal, respectively. The fifth example shows that Theorem 4 and Corollary 3 require the boundedness condition on  $i(\cdot, \cdot)$ , while the sixth shows that the same is true for Theorem 6 and Corollary 4. Example 7 shows that (9) fails if we define 1-optimality differently (see Remark 2). Finally, Example 8 shows that even under (UR), a nonstationary 1-optimal policy need not be average overtaking among stationary policies (see Remark 4).

The first two examples depend on Lemma 5 (see below), which establishes the existence of a bounded sequence of real numbers  $\{a_j\}_{j=1}^\infty$  satisfying

$$(23) \quad \begin{aligned} \liminf_N \sum_{n=1}^N \sum_{j=1}^n a_j / N &= -1, \\ \lim_{\beta \rightarrow 1^-} \sum \beta^{j-1} a_j &= \limsup_N \sum_{n=1}^N \sum_{j=1}^n a_j / N = 0. \end{aligned}$$

EXAMPLE 1. Let the state space consist of  $0, 1, 2, \dots, \infty$ . To each state there correspond two actions, 0 and 1. The functions of  $q(\cdot | \cdot, \cdot)$  and  $i(\cdot, \cdot)$  satisfy:

$$\begin{aligned} q(j + 1 | j, 0) &= q(j + 1 | j, 1) = 1, & j &= 1, 2, \dots, \\ q(1 | 0, 1) &= 1, & q(0 | 0, 0) &= q(\infty | 0, 0) = \frac{1}{2}, \\ q(\infty | \infty, 0) &= q(\infty | \infty, 1) = 1; \\ i(j, 0) &= i(j, 1) = a_{j+1}, & j &= 1, 2, \dots, \\ i(0, 1) &= a_1, & i(0, 0) &= -\frac{1}{4}, & i(\infty, 0) &= i(\infty, 1) = 0. \end{aligned}$$

Let  $\pi_j$  denote the policy which always selects action  $j$  ( $j = 0, 1$ ). One can easily show that  $\pi_1$  is optimal. That  $\pi_1$  is not average overtaking optimal follows from

$$\liminf_N \sum_{n=1}^N \sum_{j=1}^n (V^j(0, \pi_1) - V^j(0, \pi_0))/N = \liminf_N (\sum_{n=1}^N \sum_{j=1}^n a_j)/N + \frac{1}{2} \text{ and (23).}$$

EXAMPLE 2. Let the state space consist of  $0, 1, 2, \dots$ . To each state, there correspond two actions, 0 and 1. The functions  $q(\cdot | \cdot, \cdot)$  and  $i(\cdot, \cdot)$  satisfy:

$$\begin{aligned} q(1 | 0, 1) &= q(0 | 0, 0) = q(j + 1 | j, 0) = q(j + 1 | j, 1) = 1, & j &= 1, 2, \dots, \\ i(0, 1) &= a_i, & i(0, 0) &= 0, & i(j, 0) &= i(j, 1) = a_{j+1}, & j &= 1, 2, \dots. \end{aligned}$$

Define  $\pi_j$  ( $j = 0, 1$ ) as in Example 1. Using (23), one can show that  $\pi_0$  is both 1-optimal and average overtaking optimal while  $\pi_1$  is 1-optimal but not average overtaking optimal.

Our next two examples depend on the following lemma, the proof of which appears in Section 6.

LEMMA 4. *There exists a bounded sequence of real numbers  $\{v_j\}_{j=1}^\infty$  which satisfy*

$$(24) \quad \lim_N \sum_{n=1}^N \sum_{j=1}^n v_j / N = \infty, \quad \liminf_n \sum_{j=1}^n v_j / n < 0.$$

EXAMPLE 3. Let the state and action spaces be as in Example 2. Transitions are deterministic:

$$q(0|0, 0) = q(1|0, 1) = q(j+1|j, 0) = q(j+1|j, 1) = 1, \\ j = 1, 2, \dots$$

The immediate income depends only on the state:

$$i(0, 0) = i(0, 1) = 0, \quad i(j, 0) = i(j, 1) = v_j, \quad j = 1, 2, \dots$$

Define  $\pi_j$  ( $j = 0, 1$ ) as in Example 1. Clearly  $\pi_1$  is average overtaking optimal but is not lim inf average optimal. That  $\pi_1$  is optimal follows from the fact (see Hobson (1926)) that for any policy  $\pi$  and any  $s \in S$ ,

$$(25) \quad \liminf_N \sum_{n=1}^N V^n(s, \pi) / N \leq \liminf_{\beta \rightarrow 1^-} V_\beta(s, \pi) \leq \limsup_{\beta \rightarrow 1^-} V_\beta(s, \pi) \\ \leq \limsup_N \sum_{n=1}^N V^n(s, \pi) / N.$$

EXAMPLE 4. Let the state space, action space and law of transition be as in Example 3. The only difference is that we define  $i(\cdot, \cdot)$  by  $i(0, 0) = i(1, 0) = 0$  and  $i(j, 0) = i(j, 1) = w_j \equiv -v_j$  for  $j = 1, 2, \dots$ . Define  $\pi_j$  ( $j = 0, 1$ ) as in Example 1. Clearly,  $\pi_0$  is average overtaking optimal, but is not lim sup average optimal, since by (24) we have  $\lim_N \sum_{n=1}^N \sum_{j=1}^n w_j / N = -\infty$  and  $\limsup \sum_{j=1}^n w_j / n > 0$ . That  $\pi_0$  is optimal follows from (25).

EXAMPLE 5. Let the state and action space be as in Example 1. Transitions are described as follows:

$$q(\infty|\infty, 0) = q(\infty|\infty, 1) = q(0|0, 0) = 1, \quad q(1|0, 1) = \frac{1}{2} = q(\infty|0, 1), \\ q(j+1|j, 0) = q(j+1|j, 1) = q(\infty|j, 1) = q(\infty|j, 0) = \frac{1}{2}, \\ j = 1, 2, \dots$$

The immediate income depends only on the state:

$$i(0, 0) = i(0, 1) = 0 = i(\infty, 0) = i(\infty, 1), \quad i(j, 0) = i(j, 1) = 2^j w_j,$$

where  $w_j$  is defined as in Example 4 for  $j = 1, 2, \dots$ . Define  $\pi_j$  ( $j = 0, 1$ ) as in Example 1. Observe that under any stationary policy, the system eventually reaches an absorbing state with probability 1 (condition (EPR)). However, using the arguments of Example 4, one can show that  $\pi_0$  is optimal but not lim sup average optimal.

EXAMPLE 6. Let the state and action space be as in Example 1. Transitions

are described as follows:

$$\begin{aligned}
 q(\infty | \infty, 0) &= q(\infty | \infty, 1) = 1, & q(1 | 0, 1) &= q(\infty | 0, 1) = \frac{1}{2}, \\
 q(0 | 0, 0) &= q(\infty | 0, 0) = \frac{1}{2}, \\
 q(j + 1 | j, 0) &= q(j + 1 | j, 1) = q(\infty | j, 0) = q(\infty | j, 1) = \frac{1}{2}, \\
 & & j &= 1, 2, \dots
 \end{aligned}$$

The immediate income depends only on the state:

$$\begin{aligned}
 i(j, 0) &= i(j, 1) = 2^{j+1}a_{j+1}, & j &= 1, 2, \dots, \\
 i(0, 1) &= 2a_1, & i(0, 0) &= -\frac{1}{4}, & i(\infty, 0) &= i(\infty, 1) = 0,
 \end{aligned}$$

where  $a_j$ , ( $j = 1, 2, \dots$ ), satisfies (23). Define  $\pi_j$  ( $j = 0, 1$ ) as in Example 1. Using the arguments of Example 1, one can show that  $\pi_1$  is optimal, but not average overtaking optimal. Condition (UR) follows from the fact that the probability of moving to state  $\infty$  in any transition is always  $\geq \frac{1}{2}$ .

EXAMPLE 7. Let the  $S$  consist of  $0, 1, 2, \dots$  together with  $(k, 1), (k, 2), \dots, (k, 2k)$  for  $k = 1, 2, \dots$ . To each state, there correspond two actions, 0 and 1. Transitions are deterministic. When in state  $k$  ( $k = 1, 2, \dots$ ), we can either (Act 0) move to  $k + 1$  or (Act 1) move to  $(k, 1)$ . Once we reach state  $(k, 1)$ , we proceed along the path connecting  $(k, j)$  to  $(k, j + 1)$  until we hit  $(k, 2k)$ ; then we move directly to state 0. State 0 is an absorbing state. Formally,

$$\begin{aligned}
 q(0 | 0, i) &= 1 = q(0 | (k, 2k), i) = q((k, j + 1) | (k, j), i), \\
 q(k + 1 | k, 0) &= 1 = q((k, 1) | k, 1), \\
 & i = 0, 1; j = 1, \dots, 2k - 1; k = 1, 2, \dots
 \end{aligned}$$

The immediate income depends only on the state:

$$i(k, i) = 0, \quad i = 0, 1; k = 0, 1, 2, \dots,$$

and

$$\begin{aligned}
 i((k, j), i) &= 1, & 1 &\leq j \leq k, \\
 &= -2, & k < j \leq 2k, & i = 0, 1; k = 1, 2, \dots
 \end{aligned}$$

Let  $\pi_*$  denote the policy which always chooses action 0. Also, for  $N = 1, 2, \dots$ , let  $\pi_N$  be the stationary policy which selects action 1 in state  $N$  and action 0 in all other states. One can easily show that  $\lim_{\beta \rightarrow 1^-} \max_N V_\beta(1, \pi_N) > 0$ . Since  $V_\beta(1, \pi_*) = 0$  for all  $\beta$ , we have  $\liminf_{\beta \rightarrow 1^-} V_\beta(1, \pi_*) - V_{\beta^*}(1) < 0$ . Thus  $\pi_*$  is not 1-optimal if 1-optimality is defined by (13), yet  $\pi_*$  satisfies (6).

The next example depends on Lemma 5 below, the proof of which appears in Section 6.

LEMMA 5. (a) *There exists a sequence of integers  $N_0 = 0, N_1, N_2, \dots$  satisfying*

$$(26) \quad 0 < N_k^2 + 4N_k < N_{k+1}^2, \quad k = 1, 2, \dots,$$

$$(27) \quad N_k^2 \beta^{N_k} (1 - \beta)(1 - \beta^{4N_k}) \leq 2^{-k}, \quad 0 < \beta < 1; k = 1, 2, \dots,$$

and

$$(28) \quad (1 - (\frac{1}{2})^{k-1})^{I_k} \leq \frac{1}{2}, \quad k = 1, 2, \dots,$$

where  $I_k$  denotes the largest integer  $\leq (N_k^2 - N_{k-1}^2 - 4N_{k-1} - 1)/(k - 1)$ .

(b) Given  $N_0 = 0, N_1, N_2, \dots$  satisfying (26) and (27), define  $\{a_j\}_{j=1}^\infty$  by

$$(29) \quad \begin{aligned} a_j &= 0 && \text{if } N_{k-1}^2 + 4N_{k-1} < j \leq N_k^2, \\ &= -1 && \text{if } N_k^2 < j \leq N_k^2 + N_k \text{ or } N_k^2 + 3N_k < j \leq N_k^2 + 4N_k, \\ &= 1 && \text{if } N_k^2 + N_k < j \leq N_k^2 + 3N_k, \text{ for } k = 0, 1, \dots \end{aligned}$$

Then  $\{a_j\}_{j=1}^\infty$  satisfies (23).

EXAMPLE 8. Let  $N_0, N_1, N_2, \dots$  and  $a_1, a_2, \dots$  satisfy the conditions of Lemma 5, and let  $\alpha_1, \alpha_2, \dots$  satisfy

$$(30) \quad \alpha_k N_k^2 \geq \frac{1}{2}, \quad 0 < \alpha_k < 1, \quad k = 1, 2, \dots$$

Let the state space consist of  $1, 2, \dots$  together with  $(k, 1), \dots, (k, 4N_k)$  for  $(k = 1, 2, \dots)$ . To each state, there correspond four actions: 0, 1, 2 and 3. When in state  $k$ , we can either (Act 0) move directly to state  $(k, 1)$ , (Act 1) move directly to state 1, (Act 2) move to states  $k + 1$  and 1 with probability  $\frac{1}{2}$  each by tossing a fair coin, or (Act 3) remain in state  $k$  with probability  $\alpha_k$  and move to state 1 with probability  $1 - \alpha_k$  by tossing a biased coin. Specifically,

$$\begin{aligned} q((k, 1) | k, 0) &= 1 = q(1 | k, 1), & q(k + 1 | k, 2) &= \frac{1}{2} = q(1 | k, 2), \\ q(k | k, 3) &= \alpha_k = 1 - q(1 | k, 3), & & k = 1, 2, \dots \end{aligned}$$

Once we reach state  $(k, 1)$ , we proceed along the path connecting  $(k, j)$  to  $(k, j + 1)$  until we hit  $(k, 4N_k)$ ; then we move directly to state 1, i.e.,

$$\begin{aligned} q((k, j + 1) | (k, j), i) &= q(1 | (k, 4N_k), i) = 1, \\ i &= 0, \dots, 3; j = 1, \dots, 4N_k - 1; k = 1, 2, \dots \end{aligned}$$

All coin tosses are assumed independent. The immediate income depends only on the state:

$$\begin{aligned} i(k, i) &= 0, & i &= 0, 1, 2; k = 1, 2, \dots, \\ i((k, j), i) &= -1, & 1 &\leq j \leq N_k \text{ or } 3N_k < j \leq 4N_k \\ &= 1, & N_k &< j \leq 3N_k, \quad i = 0, 1; k = 1, 2, \dots \end{aligned}$$

One can show that given any (nonrandomized) stationary policy, the state space of the associated Markov chain can be expressed as a union of uniformly recurrent classes (see Section 2). These classes are, of course, not disjoint. In fact, for each such class, the set of positive recurrent states consists of the states accessible from state 1. The role played by  $j$  in the definition of uniform recurrence (see Section 2) is assumed by state 1. Hence condition (UR) holds.

Let  $\pi_1$  denote the policy which always selects action 1. One can easily show that  $\pi_1$  is both optimal and average overtaking optimal. We shall construct a

nonstationary policy  $\pi_*$  satisfying

$$(31) \quad \sum_{j=1}^{\infty} \beta^{j-1} a_j \leq \sum_{j=1}^{\infty} \beta^{j-1} r_j(1, \pi_*) \leq 0, \quad 0 \leq \beta < 1,$$

$$(32) \quad \sum_{n=1}^N V^n(1, \pi_*) \leq \left(\frac{1}{4}\right) \sum_{n=1}^N \left(\sum_{j=1}^n a_j\right), \quad N = 1, 2, \dots,$$

and

$$(33) \quad r_j(s, \pi_*) = r_j(s, \pi_1), \quad s \neq 1; j = 1, 2, \dots$$

Now, (23), (31) and (32) imply

$$(34) \quad \liminf_N \sum_{n=1}^N V^n(1, \pi_*)/N \leq -\frac{1}{4} \quad \text{and} \\ \lim_{\beta \rightarrow 1^-} \sum_{j=1}^n \beta^{j-1} r_j(1, \pi_*) = 0.$$

Using (33), (34), and the fact that  $\pi_1$  is both optimal and average overtaking optimal, one can show that any such  $\pi_*$  is 1-optimal but not average overtaking optimal. We leave the details to the reader.

Now we construct  $\pi_*$ . Given any initial state  $s \neq 1$ , always choose *act 1*. (Clearly, (33) holds.) For the case in which the initial state is 1, follow the program:

I. *Select act 3 from time 1 until time  $N_1^2$ . Select act 0 from time  $N_1^2$  until time  $N_1^2 + 4N_1 + 1$ . Then set  $K = 1$ .*

II. *Select act 2 from time  $N_K^2 + 4N_K + 1$  until time  $N_{K+1}$  or the first time that the system reaches state  $K + 1$ , whichever occurs first.*

III. *If the system reaches state  $K + 1$  before time  $N_{K+1}^2$ , then choose act 3 from the time the system reaches  $K + 1$  until time  $N_{K+1}^2$ .*

IV. *If the system occupies state  $K + 1$  at time  $N_{K+1}^2$ , then choose act 0 from time  $N_{K+1}^2$  until time  $N_{K+1}^2 + 4N_{K+1} + 1$ ; otherwise, choose act 1 from time  $N_{K+1}^2$  until time  $N_{K+1}^2 + 4N_{K+1} + 1$ .*

V. *Set  $K = K + 1$ . Then go to II.*

We need only establish (31) and (32). Assume  $\pi_*$  is used and the initial state is 1. For any  $k$  ( $k = 1, 2, \dots$ ), one can easily show that if the system occupies state  $k$  at time  $N_k^2$ , then the income for period  $j$  will equal  $a_j$  for values of  $j \in [N_{k-1}^2 + 4N_{k-1} + 1, N_k^2 + 4N_k]$ ; while if the system does not occupy state  $k$  at time  $N_k^2$ , then the income will be zero for all such values of  $j$ . Using this and (29), one can establish (31). As the reader can verify, (32) will also follow once we establish

$$(35) \quad P[\text{system occupies state } k \text{ at time } N_k^2] \geq \frac{1}{4}, \quad k = 1, 2, \dots$$

To prove (35), we need

$$(36) \quad P[\text{system reaches state } k \text{ by time } N_k^2] \geq \frac{1}{2}, \quad k = 1, 2, \dots$$

Toward establishing (36), observe that each time act 2 is chosen a fair coin is tossed. Call a toss “favorable” if it does not direct the system to state 1. Now, one can show that the system will always be in state 1 at time  $N_{k-1}^2 + 4N_{k-1} + 1$ , ( $k = 1, 2, \dots$ ). Hence one way in which the system can reach state  $k$  by time

$N_k^2$  is for one of the first  $I_k$  groups of successive coin tosses to constitute a favorable run of length  $k - 1$  ( $I_k$  is defined in Lemma 5(a)). The probability of none of those runs being favorable equals  $(1 - (\frac{1}{2})^{k-1})^{I_k}$  by the independence of coin tosses. Applying (28), we get (36). Now, given that the system reaches state  $k$  at time  $\tau < N_k^2$ , the system remains in state  $k$  between time  $\tau$  and time  $N_k^2$  with probability  $\alpha_k^{N_k^2 - \tau} > \alpha_k^{N_k^2}$ , which is  $\geq \frac{1}{2}$  by (28). This and (36) give us (35).

**6. Proofs.** The only things left are the proofs of Lemmas 1, 4 and 5.

**PROOF OF LEMMA 1.** Let  $\{a_j\}_{j=1}^\infty$  be any sequence satisfying (16). It is convenient to work with the sequence  $\{x_j\}_{j=1}^\infty$  defined by

$$(37) \quad x_j = a_j - \liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{i=1}^\infty \beta^{i-1} a_i, \quad j = 1, 2, \dots$$

Clearly,

$$(38) \quad \liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{j=1}^\infty \beta^{j-1} x_j = 0, \quad |x_j| \leq 2M, \quad j = 1, 2, \dots$$

By (37) and (38), (15) is equivalent to

$$(39) \quad \liminf \sum_{j=1}^n x_j/n > -\epsilon,$$

while (17) is equivalent to

$$(40) \quad \limsup \sum_{j=1}^n x_j/n < \delta.$$

These equivalences allow us to establish the lemma by showing that

$$(41) \quad \begin{aligned} \liminf_n \sum_{j=1}^n x_j/n &\leq -\epsilon, \\ \limsup_n \sum_{j=1}^n x_j/n &< \delta(\epsilon, M) \equiv \epsilon(1 - e^{-\epsilon/3M})/6\epsilon \end{aligned}$$

is impossible, since it will then follow that  $\delta = \delta(\epsilon, M)$  meets the required conditions. (The reader should verify this.)

Assume (41). Let  $S_n = \sum_{j=1}^n x_j$  ( $n = 1, 2, \dots$ ) and let  $\rho = \limsup S_n/n$ . Clearly, there exists an integer  $N$  such that

$$(42) \quad S_n \leq 2\rho n, \quad n \geq N.$$

By (41), there exist integers  $N_1, N_2, \dots$  satisfying

$$(43) \quad N < N_j < N_{j+1}, \quad \lim_j N_j = \infty, \quad S_{N_j} \leq -(\frac{2}{3})\epsilon N_j, \quad j = 1, 2, \dots$$

Now, (43) and the fact that the  $x_j$ 's are bounded in absolute value by  $2M$  (see (38)) imply that

$$(44) \quad S_n \leq -(\frac{1}{3})\epsilon N_j, \quad L_j \leq n \leq U_j,$$

where  $L_j$  is the smallest integer  $\geq N_j(1 - \epsilon/6M)$  and  $U_j$  is the largest integer  $\leq N_j(1 + \epsilon/6M)$ , ( $j = 1, 2, \dots$ ). It is easy to show that

$$(1 - \beta) \sum_{j=1}^\infty \beta^{j-1} x_j = (1 - \beta)^2 \sum_{j=1}^\infty \beta^{j-1} S_j.$$

Clearly, the RHS equals

$$(1 - \beta)^2 (\sum_{j=1}^{N_1} \beta^{j-1} S_j + \sum_{j=L_2}^{N_2} \beta^{j-1} S_j + \sum_{j=L_k}^{U_k} \beta^{j-1} S_j + \sum_{j=U_{k+1}}^\infty \beta^{j-1} S_j),$$

which according to (38), (42) and (44) is bounded above by

$$2(1 - \beta)^2 N^2 M + (1 - \beta)^2 \sum_{j=1}^{\infty} \beta^{j-1} (2\rho j) - (1 - \beta)^2 \sum_{j=L_k+1}^{U_k} \beta^{j-1} (\frac{1}{3})^\epsilon N_k$$

for  $k = 1, 2, \dots$ . Consequently,

$$(45) \quad (1 - \beta) \sum_{j=1}^{\infty} \beta^{j-1} x_j \leq 2(1 - \beta)^2 N^2 M + 2\rho - (\frac{1}{3})^\epsilon N_k (1 - \beta) \beta^{L_k} (1 - \beta^{U_k - L_k}),$$

for  $k = 1, 2, \dots$ . Let  $\beta_k = 1 - 1/N_k$ , ( $k = 1, 2, \dots$ ). Evidently,  $N_k(1 - \beta_k) = 1$  and  $\lim_k \beta_k^{N_k} = e^{-1}$ . Now, one can easily show that

$$\begin{aligned} \lim_k N_k (1 - \beta_k) \beta_k^{L_k} (1 - \beta_k^{U_k - L_k}) &= \lim_k \beta_k^{L_k} - \beta_k^{U_k} = e^{-(1-\epsilon/6M)} - e^{-(1+\epsilon/6M)} > (1 - e^{-\epsilon/3M})/e. \end{aligned}$$

Thus, the fact that  $\lim_k (1 - \beta_k)^2 N^2 M = 0$ , and (45) imply

$$\limsup_k (1 - \beta_k) \sum_{j=1}^{\infty} (\beta_k)^{j-1} x_j < 2\rho - \epsilon(1 - e^{-\epsilon/3M})/3e.$$

Consequently,  $\liminf_{\beta \rightarrow 1^-} (1 - \beta) \sum_{j=1}^{\infty} \beta^{j-1} x_j < 0$  whenever  $\rho < \delta(\epsilon, M)$  (see (41)), contradicting (38). Lemma 1 follows.

**PROOF OF LEMMA 4.** Our proof is constructive. For each  $k$  ( $k = 0, 1, \dots$ ), define  $r_j$  by

$$\begin{aligned} r_j &= 1, & 2(2^k - 1) < j \leq 2(2^k - 1) + 2^k \\ &= -1, & 2(2^k - 1) + 2^k < j \leq 2(2^{k+1} - 1). \end{aligned}$$

Let  $v_j(\delta) = r_j - \delta$  for all real  $\delta$ . Clearly,  $\liminf_n \sum_{j=1}^n v_j(\delta)/n = -\delta$ . Hence, all we need is a  $\delta > 0$  satisfying

$$(46) \quad \liminf_N \sum_{n=1}^N \sum_{j=1}^n v_j(\delta)/N = \infty.$$

We will show that  $\delta = 2^{-11}$  works. Let  $S_n = \sum_{j=1}^n r_j$  ( $n = 1, 2, \dots$ ) and  $m(k) = 2(2^k - 1)$ , ( $k = 0, 1, \dots$ ). For any  $\delta > 0$ ,

$$(47) \quad \sum_{j=1}^N (S_j - j\delta)/N \geq \sum_{j=1}^{m(k)} S_j/m(k+1) - \sum_{j=1}^{m(k+1)} j\delta/m(k),$$

$$m(k) \leq N \leq m(k+1).$$

Now, the RHS of (47) equals

$$(\sum_{j=1}^{m(k)} S_j/m(k))(m(k)/m(k+1)) - \sum_{j=1}^{m(k)} j\delta/m(k)(\sum_{j=1}^{m(k+1)} j/\sum_{j=1}^{m(k)} j).$$

Using the well-known formula for the sum of consecutive positive integers, one gets  $\sum_{j=1}^{m(k)} j = (2^k - 1)(2^{k+1} - 1)$ , ( $k = 1, 2, 3, \dots$ ). Also,  $2^{k-1} \leq 2^k - 1$ , ( $k = 1, 2, \dots$ ). Hence,  $\sum_{j=1}^{m(k+1)} j/\sum_{j=1}^{m(k)} j \leq 8$  and  $m(k)/m(k+1) \geq \frac{1}{4}$ , ( $k = 1, 2, \dots$ ). This implies that the RHS of (47) is  $\geq (\frac{1}{4}) \sum_{j=1}^{m(k)} (S_j - 32j\delta)/m(k)$ , ( $k = 1, 2, \dots$ ). Consequently, if  $\delta > 0$  satisfies

$$(48) \quad \lim_k \sum_{j=1}^{m(k)} (S_j - j\delta)/m(k) = \infty,$$

then  $2^{-5\delta}$  satisfies (46). We will show that  $\delta = 2^{-6}$  satisfies (48).

As the reader can easily verify,

$$\begin{aligned} S_{m(j)+h} &= h, & \text{if } h &= 1, \dots, 2^j, \\ &= 2^{j+1} - h, & \text{if } h &= 2^j + 1, \dots, 2^{j+1}, \end{aligned}$$

( $j = 0, 1, \dots$ ). This gives us

$$\sum_{i=m(k-1)+1}^{m(k)} S_i = 2^{2(k-1)}, \quad k = 1, 2, \dots$$

Using the fact that  $\sum_{j=1}^{m(k)} j \leq 2^{2k+1}$ , we have

$$(49) \quad \sum_{j=1}^{m(k)} (S_j - j\delta) \geq 2^{2k-4} - 2^{2k+1}\delta$$

whenever  $\delta > 0$ . Observe that for  $\delta = 2^{-6}$ , the RHS of (49)  $\geq 2^{2k-5}$ . Consequently, for  $\delta = 2^{-6}$ ,

$$\liminf_k \sum_{j=1}^{m(k)} (S_j - j\delta)/m(k) \geq \liminf_k (2^{2k-5})2^{-k-1} = \infty.$$

Lemma 5 requires the following result:

LEMMA 6. For any  $\varepsilon > 0$ , there exists a positive integer  $N_0$  such that  $N \geq N_0$  implies

$$N^2\beta^{N^2}(1 - \beta)(1 - \beta^{4N}) \leq \varepsilon, \quad 0 < \beta < 1.$$

PROOF. Let  $\rho = -\ln \beta$ . We need only show that

$$F(\rho, x) = x^2e^{-\rho x^2}(1 - e^{-\rho})(1 - e^{-4\rho x})$$

converges to 0 uniformly in  $\rho$ , ( $0 \leq \rho < \infty$ ), as  $x$  approaches  $\infty$ . We will do this by proving that the existence of a  $\delta > 0$  and sequences  $\{x_j\}_{j=1}^\infty$  and  $\{\rho_j\}_{j=1}^\infty$  satisfying

$$(50) \quad \lim_j x_j = \infty, \quad x_j < x_{j+1}, \quad F(x_j, \rho_j) \geq \delta, \quad j = 1, 2, \dots$$

leads to a contradiction. To begin with, whenever  $\{x_j\}_{j=1}^\infty$  and  $\{\rho_j\}_{j=1}^\infty$  satisfy (50), we have

$$(51) \quad \lim_j \rho_j x_j = 0, \quad \lim_j \rho_j = 0.$$

Otherwise, there exist a  $\lambda > 0$  and subsequences  $\{x_{n_j}\}_{j=1}^\infty$  and  $\{\rho_{n_j}\}_{j=1}^\infty$  such that  $x_{n_j}\rho_{n_j} \geq \lambda$ , ( $j = 1, 2, \dots$ ). But then  $F(\rho_{n_j}, x_{n_j}) \leq x_{n_j}^2 e^{-\lambda x_{n_j}}$ , ( $j = 1, 2, \dots$ ), which is impossible by (50). Now (51) and the fact that  $1 - e^{-z} \leq ze^z$  for all positive  $z$  imply that for  $j$  sufficiently large,

$$F(\rho_j, x_j) \leq (5\rho_j x_j) \cdot (\rho_j x_j^2 e^{-\rho_j x_j^2}).$$

But the second factor in the RHS is bounded while the first converges to 0 by (51). This contradicts (50), finishing the proof.

PROOF OF LEMMA 5. First, part (a). Evidently there exist integers  $I(1), I(2), \dots$  such that

$$(52) \quad (1 - (\frac{1}{2})^{k-1})^I \leq \frac{1}{2}, \quad I \geq I(k), \quad k = 2, 3, \dots$$

By Lemma 6, there exist integers  $N_1, N_2, \dots$  satisfying (27). Lemma 6 also allows us to choose this sequence so that  $N_k^2 \geq (k - 1)I(k) + N_{k-1}^2 + 4N_{k-1} + 1$  for  $k = 1, 2, \dots$ , where  $I(k)$  satisfies (52) and  $N_0 = 0$ . One can easily show that any such sequence must satisfy (26) and (28). Part (a) follows.



For part (b), let  $N_0 = 0, N_1, N_2, \dots$  satisfy (26) and (27), and let  $J_k = \{j | N_k^2 + 1 \leq j \leq N_k^2 + 4N_k\}$  for  $k = 1, 2, \dots$ . Define  $s_i = \sum_{j=1}^i a_j$  and  $T_i = \sum_{j=1}^i s_j, (i = 1, 2, \dots)$ . Evidently, if  $i \in J_k$  for some  $k$ , then  $s_i$  satisfies

$$(53) \quad \begin{aligned} s_{N_k^2+j} &= -j && \text{if } 0 \leq j \leq N_k \\ &= -(2N_k - j) && \text{if } N_k \leq j \leq 3N_k \\ &= 4N_k - j && \text{if } 3N_k \leq j \leq 4N_k; \end{aligned}$$

while if  $i \notin \bigcup_k J_k$ , then  $s_i = 0$ . Using (53) and the well-known formula for the sums of consecutive positive integers, one can show that  $T_{N_k^2+2N_k} = -N_k^2$  and  $-N_k^2 \leq T_j \leq 0, (j \in J_k)$ , for  $k = 1, 2, \dots$ . Also,  $T_j = 0, j \notin \bigcup_k J_k$ . One can easily show that these facts imply  $\limsup_N T_N/N = 0$  and  $\liminf_N T_N/N = \lim_k (T_{N_k^2+2N_k})/(N_k^2 + 2N_k) = -1$ . The only thing left to establish is

$$(54) \quad \lim_{\beta \rightarrow 1^-} \sum_{j=1}^{\infty} \beta^{j-1} a_j = 0.$$

We will make use of the identity

$$(55) \quad (1 - \beta)^2 \sum_{j=1}^{\infty} \beta^{j-1} T_j = \sum_{j=1}^{\infty} \beta^{j-1} a_j, \quad 0 < \beta < 1.$$

Since  $T_j \geq -N_k^2$  for  $j \in J_k$ , we have  $\sum_{j \in J_k} \beta^{j-1} T_j \geq -N_k^2 \beta^{N_k^2} (\sum_{j=1}^{4N_k} \beta^{j-1})$ . Hence,

$$(56) \quad (1 - \beta)^2 \sum_{j \in J_k} \beta^{j-1} T_j \geq -N_k^2 \beta^{N_k^2} (1 - \beta)(1 - \beta^{4N_k}).$$

Let  $\delta > 0$  be arbitrary. Condition (54) will follow if we can select a  $\beta_0 < 1$  so that

$$(57) \quad \sum_{j=1}^{\infty} \beta^{j-1} a_j \geq -\delta, \quad \beta_0 \leq \beta < 1.$$

Select  $k_0$  and  $\beta_0 < 1$  so that  $\sum_{k=k_0+1}^{\infty} 2^{-k} \leq \delta/2$  and  $(1 - \beta)N_{k_0}^2(1 - \beta^{N_{k_0}^2}) \leq \delta/2$  for  $\beta_0 \leq \beta < 1$ . The identity (55) implies that

$$(58) \quad \sum_{j=1}^{\infty} \beta^{j-1} a_j = (1 - \beta)^2 \sum_{j=1}^{N_{k_0}^2} \beta^{j-1} T_j + (1 - \beta)^2 \sum_{j=N_{k_0}^2+1}^{\infty} \beta^{j-1} T_j.$$

By (56), the fact that  $T_j = 0$  for  $j \notin \bigcup_k J_k$ , and the fact that  $t_j \geq -N_k^2$  for  $j \in J_k$  ( $k = 1, 2, \dots$ ), the RHS of (58) is  $\geq (1 - \beta)^2 \sum_{j=1}^{N_{k_0}^2} N_{k_0}^2 \beta^{j-1} - \sum_{k=k_0+1}^{\infty} N_k^2 \beta^{N_k^2} (1 - \beta)(1 - \beta^{4N_k})$ . But the latter is  $\geq -(1 - \beta)N_{k_0}^2(1 - \beta^{N_{k_0}^2}) - \delta/2$  by (27) and our choice of  $k_0$ . This and our method of selecting  $\beta_0$  give us (57). We are done.

**Acknowledgments.** We wish to thank Ben Fox for his helpful comments.

REFERENCES

BLACKWELL, D. (1962). Discrete dynamic programming. *Ann. Math. Statist.* **33** 719-726.  
 BLACKWELL, D. (1965). Discounted dynamic programming. *Ann. Math. Statist.* **36** 226-235.  
 BROWN, B. (1965). On the iterative method of dynamic programming on a finite space discrete time Markov process. *Ann. Math. Statist.* **36** 1279-1285.  
 CHUNG, K. L. (1967). *Markov Chains with Stationary Transition Probabilities*. Springer-Verlag, Berlin.  
 DENARDO, E. and MILLER, B. (1968). An optimality condition for discrete dynamic programming with no discounting. *Ann. Math. Statist.* **39** 1220-1227.  
 DERMAN, C. (1964). On sequential control processes. *Ann. Math. Statist.* **35** 341-349.

- DERMAN, C. (1966). Denumerable state Markovian decision processes—average cost criterion. *Ann. Math. Statist.* **37** 1545–1553.
- DERMAN, C. and VEINOTT, A. (1967). A solution to a countable system of equations arising in Markovian decision processes. *Ann. Math. Statist.* **38** 582–585.
- DOOB, J. (1953). *Stochastic Processes*. Wiley, New York.
- FLYNN, J. (1974). Averaging versus discounting in dynamic programming: a counterexample. *Ann. Statist.* **2** 411–413.
- HOBSON, E. (1926). *The Theory of Functions of a Real Variable and the Theory of Fourier's Series*. Cambridge Univ. Press.
- LIGGETT, T. and LIPPMAN, S. (1969). Stochastic games with perfect information and time average payoff. *SIAM Rev.* **11** 604–607.
- LIPPMAN, S. (1969). Criterion equivalence in discrete dynamic programming. *Operations Res.* **17** 920–923.
- MAITA, A. (1965). Dynamic programming for countable state systems. *Sankhyā Ser. A* **27** 259–266.
- OREY, S. (1971). *Limit Theorems for Markov Chain Transition Probabilities*. Van Nostrand, Princeton.
- ROSS, S. (1968). Non-discounted denumerable Markovian decision models. *Ann. Math. Statist.* **39** 412–423.
- ROSS, S. (1971). On the nonexistence of  $\epsilon$ -optimal randomized stationary policies in average cost Markov decision models. *Ann. Math. Statist.* **42** 1567–1568.
- STRAUCH, R. (1966). Negative dynamic programming. *Ann. Math. Statist.* **37** 871–890.
- VEINOTT, A. (1966). On finding optimal policies in discrete dynamic programming with no discounting. *Ann. Math. Statist.* **37** 1284–1294.
- VEINOTT, A. (1969). Discrete dynamic programming with sensitive discount optimality criteria. *Ann. Math. Statist.* **40** 1635–1660.

SCHOOL OF BUSINESS ADMINISTRATION  
WAYNE STATE UNIVERSITY  
DETROIT, MICHIGAN 48202