



---

On Optimal Policies and Martingales in Dynamic Programming

Author(s): Ulrich Rieder

Source: *Journal of Applied Probability*, Vol. 13, No. 3 (Sep., 1976), pp. 507-518

Published by: Applied Probability Trust

Stable URL: <https://www.jstor.org/stable/3212470>

Accessed: 26-10-2019 19:39 UTC

---

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



*Applied Probability Trust* is collaborating with JSTOR to digitize, preserve and extend access to *Journal of Applied Probability*

JSTOR

## ON OPTIMAL POLICIES AND MARTINGALES IN DYNAMIC PROGRAMMING

ULRICH RIEDER, *University of Hamburg*

### Abstract

A martingale approach to a dynamic program with general state and action spaces is taken. Several necessary and sufficient conditions are given for a policy to be optimal. The results comprehend and modify different criteria of optimality given for dynamic programming problems. Finally, two applications are stated.

DYNAMIC PROGRAMMING; GAMBLING; EXCESSIVE FUNCTION; MARTINGALE; OPTIMAL POLICY; THRIFTY POLICY; EQUALIZING POLICY

### 1. Introduction and summary

Dubins and Savage (1965) and Sudderth (1972) have given an elegant characterization of optimal strategies for general gambling models. In the present paper we take a martingale approach to dynamic programming problems and derive similar characterizations of optimal policies in dynamic programs with general state and action spaces. The impetus for our martingale approach came from Blackwell (1970) and Hordijk (1974).

The analysis is based only on methods of martingale theory. However, most of the results can also be proved by dynamic programming methods (cf. Rieder (1975)).

In Section 2 we formulate the dynamic programming model. Assumptions are stated that ensure the existence of the conditional expected rewards; in that connection we remark that the reward function is not assumed to be bounded from above or from below. Section 3 deals with excessive functions relative to the dynamic program. Some examples and properties of excessive functions are listed and an optional sampling theorem for decision processes is proved. In Section 4 we give several necessary and sufficient conditions for a general (not necessarily Markov) policy to be  $(s, \varepsilon)$ -optimal (resp.  $s$ -optimal), e.g. it is shown that a policy is  $s$ -optimal if and only if it is  $s$ -thrifty and  $s$ -equalizing. Thrifty policies and equalizing policies are characterized in Section 5 and Section 6, respectively. The results are closely related to results in Dubins and Savage

---

Received in revised form 14 January 1976.

(1965) and Sudderth (1972). They comprehend and modify different criteria of optimality given by Blackwell ((1965), Theorem 6f), Hinderer ((1970), Theorems 17.1, 17.6 and 17.7), Chow, Robbins and Siegmund ((1971), Theorem 4.10), Hordijk ((1974), Theorems 4.6, 6.3 and 6.5), Ross ((1974), Propositions 1 and 2) and Schäl ((1975), Theorems 5.2 and 5.3). In particular, the results are true both in the positive and negative case. The final Section 7 contains two applications of our new results. One is to generalize a theorem of Blackwell (1970), the other is to give rather weak sufficient conditions guaranteeing the existence of an optimal policy. We remark that our methods can also be used to derive similar characterizations for non-stationary non-Markovian dynamic programs.

## 2. The dynamic programming model

In the present paper we consider a stationary dynamic program  $((S, \mathcal{S}), (A, \mathcal{A}), D, q, r)$  of the following meaning:

(i) The *state space*  $(S, \mathcal{S})$  and the *action space*  $(A, \mathcal{A})$  are assumed to be standard Borel spaces.

(ii) The set  $D$  is a measurable subset of  $S \times A$  which contains the graph of a measurable map from  $S$  into  $A$ . For any  $s \in S$  the non-empty and measurable  $s$ -section  $D(s)$  of  $D$  is called the *set of all admissible actions* if the system is in state  $s$ .

(iii)  $q$  is a transition probability from  $D$  to  $S$ , the so-called *transition law* during a single stage.

(iv) The *reward function*  $r$  is a measurable map from  $D$  into the set  $\mathbb{R}$  of real numbers.

*Remark.* In this paper we do not consider any discount factor. If the reward function for the time period  $(n, n + 1]$  has the form

$$\beta(s_1, a_1, s_2) \cdots \beta(s_{n-1}, a_{n-1}, s_n) r(s_n, a_n)$$

where  $\beta$  is a bounded measurable function from  $D \times S$  into the set of the non-negative real numbers and may be interpreted as a (generalized) discount factor (cf. Schäl (1975)), the results of this paper can be easily extended to this more general case.

We write  $H_1 = S$ ,  $H_{n+1} = D \times H_n$ ,  $n \in \mathbb{N}$ .  $\mathbb{N}$  is the set of the positive integers. As usual, a *randomized policy*  $\pi = (\pi_n)$  is defined as a sequence of transition probabilities  $\pi_n$  from  $H_n$  to  $A$  such that  $\pi_n(h, D(s_n)) = 1$  for all  $h = (s_1, a_1, s_2, \dots, s_n) \in H_n$ ,  $n \in \mathbb{N}$ . Let  $\Delta$  denote the set of all randomized policies. We write  $F$  for the set of all *decision functions*, i.e. of all measurable maps  $f: S \rightarrow A$  whose graphs belong to  $D$ . Then a (deterministic) *Markov policy* is a sequence  $(f_n)$  with  $f_n \in F$ ,  $n \in \mathbb{N}$ . A (deterministic) *stationary policy* is a Markov policy  $(f_n)$  where  $f_n = f$  is independent of  $n$ . For such a policy we write  $f^\infty$ .

Given some initial state  $s$ , any policy  $\pi$  defines a probability measure  $P_\pi$  on the product space  $H = S \times A \times S \times A \cdots$  endowed with the product  $\sigma$ -algebra  $\mathcal{H} = \mathcal{S} \otimes \mathcal{A} \otimes \mathcal{S} \otimes \mathcal{A} \cdots$  and thus a random process  $(\zeta_1, \alpha_1, \zeta_2, \alpha_2, \cdots)$  (cf. Hinderer (1970), p. 80) where  $\zeta_n$  and  $\alpha_n$  denote the projection from  $H$  onto the  $n$ th state space and the  $n$ th action space, respectively. Then  $\chi_n = (\zeta_1, \alpha_1, \zeta_2, \cdots, \zeta_n)$  describes the history at time  $n$ .

In order that the total expected rewards are well-defined we have to impose some convergence assumptions. Let  $i \in \{-, +\}$ . Define

$$W_i(s) = \sup_{\pi \in \Delta} E_\pi \left[ \sum_1^\infty r_n^i \right] \quad s \in S,$$

where  $r_n = r \circ (\zeta_n, \alpha_n)$  and where the expectation is taken with respect to  $P_\pi$ ,  $c^+ = \max(0, c)$  for any extended real number  $c$ . We shall make use of the following assumption.

*Assumption  $A^+$ .*  $W_i(s) < \infty$ ,  $s \in S$ .

Throughout the paper we impose a general assumption.

*General Assumption.* Either  $A^+$  or  $A^-$  holds.

Our general assumption implies that the total reward  $\sum_1^\infty r_n$  is  $P_\pi$ -quasi-integrable for any policy  $\pi$  and  $s \in S$  and therefore the following functions

$$V_\pi(s) = E_\pi \left[ \sum_1^\infty r_n \right] \quad \pi \in \Delta, \quad s \in S$$

are well-defined.

Formulations simplify considerably by the use of the following isotone operators. Denote by  $M_i$ ,  $i \in \{-, +\}$  the set of all measurable functions  $u$  from  $S$  into the set  $\bar{\mathbb{R}}$  of the extended real numbers such that  $u^+ < \infty$  and  $E_\pi u^+ \circ \zeta_n < \infty$ ,  $n \in \mathbb{N}$ ,  $s \in S$ ,  $\pi \in \Delta$ . If Assumption  $A^+$  holds, the operators  $L_f$ ,  $f \in F$  and  $U$  are well-defined on  $M_i$  by

$$L_f u(s) = r(s, f(s)) + \int q(s, f(s), ds') u(s')$$

$$Uu(s) = \sup_{f \in F} L_f u(s), \quad s \in S.$$

### 3. Excessive functions and martingales

A function  $v \in M_i$  is called *excessive* for the dynamic program if

$$(3.1) \quad v \geq V_\pi, \quad \pi \in \Delta,$$

$$(3.2) \quad v \geq Uv.$$

We use  $\mathcal{E}_i$  to denote the set of all excessive functions. The definition above is a slight modification of the concept introduced by Hordijk (1974); in particular Hordijk assumes the reward structure to be a charge structure and considers only Markov policies. Strauch (1966) has used another definition of an excessive function.

A measurable map  $\tau: H \rightarrow \mathbb{N} + \{\infty\}$  is called a *stopping time* relative to the sequence of  $\sigma$ -algebras  $\mathcal{F}_n = \sigma(\chi_n)$  if  $[\tau = n] \in \mathcal{F}_n$  for all  $n \in \mathbb{N}$ . We write  $T$  for the set of all stopping times. For later use we set  $\mathcal{F}_\infty = \sigma(\zeta_1, \alpha_1, \zeta_2, \dots)$ .

Throughout this paper we shall assume that the set  $\mathcal{E}_i$  is not empty. This assumption holds in many practical problems. Let us consider three *examples*.

(i) The so-called *value-function*  $V(s) = \sup_{\pi \in \Delta} V_\pi(s)$ ,  $s \in S$ , need not belong to  $\mathcal{E}_i$ . But, if  $V$  is measurable then  $V \in \mathcal{E}_i$  (cf. Hinderer (1971)) and  $V$  is the smallest function in  $\mathcal{E}_i$ , i.e.  $V \leq v$  for all  $v \in \mathcal{E}_i$ . The function  $V$  is measurable, e.g., if there exists a policy  $\pi$  such that  $V_\pi = V$ . We recall that  $V$  is always universally measurable (cf. Hinderer (1971)).

(ii) The *limit function*  $V_\infty = \liminf_n V_n$ , where  $V_n(s) = \sup_{\pi \in \Delta} E_{\pi s} [\sum_{k=1}^n r_k]$ ,  $n \in \mathbb{N}$ ,  $s \in S$ , satisfies the inequality  $V_\infty \geq V_\pi$  for all  $\pi$ . Thus, if  $V_\infty$  is measurable and if

$$(3.3) \quad \liminf_n L_f V_n \geq L_f V_\infty \quad \text{for all } f \in F$$

holds, then  $V_\infty \in \mathcal{E}_i$ . (For a proof of (3.2) use the relation  $V_{n+1} \geq L_f V_n$ ,  $f \in F$ .) Under Assumption  $A^-$  the Relation (3.3) is satisfied by Fatou's Lemma.

(iii) Let  $g$  be a measurable map from  $S$  into  $\mathbb{R}$ . Then  $V^*: S \rightarrow \mathbb{R}$

$$V^*(s) = \sup_{\pi \in \Delta} \sup_{\tau \in T} E_{\pi s} \left[ \sum_1^{\tau-1} r_n + (g \circ \zeta_\tau 1_N \circ \tau) \right]^+, \quad s \in S$$

belongs to  $\mathcal{E}_i$ , provided that the expectations on the right exist and that  $V \in M_i$ .

Now we list some *properties* of excessive functions.

**Remark 3.1.** If  $v$  is the smallest excessive function in  $\mathcal{E}_i$  and  $Uv \in M_i$ , then  $v = Uv$ .

**Remark 3.2.** If  $v \in \mathcal{E}_i$  and  $v = L_f v$  for some  $f \in F$ , then  $v = Uv$ .

**Remark 3.3.** If  $v \in \mathcal{E}_i$  then  $\liminf_n E_{\pi s} v \circ \zeta_n \geq 0$  on  $\{s: V_\pi(s) \in \mathbb{R}\}$ ,  $\pi \in \Delta$ .

**Remark 3.4.** If  $A^-$  holds and if  $v \in \mathcal{E}_i$  then  $\lim_n E_{\pi s} v \circ \zeta_n = 0$ ,  $s \in S$ ,  $\pi \in \Delta$ .

**Remark 3.5.** (cf. Hordijk (1974), Theorem 2.17.) Let  $v \in M_i$  such that  $v(s) < \infty$ ,  $s \in S$ . Then  $v \in \mathcal{E}_i$  iff  $v \geq Uv$  and  $\lim_n E_{\pi s} v \circ \zeta_n \geq 0$  on  $\{s: V_\pi(s) > -\infty\}$  for all  $\pi$ .

<sup>†</sup> Let  $1_B$  denote the indicator function of the set  $B$ .

For the remainder of this paper let  $v \in \mathcal{E}_i$  be fixed and suppose that Assumption  $A'$  holds. Then the random variables  $X_n$

$$X_n = \sum_{k=1}^{n-1} r_k + v \circ \zeta_n, \quad n \in \mathbb{N}, \quad X_\infty = \sum_{k=1}^{\infty} r_k$$

are  $P_{\pi_s}$ -quasi-integrable for any  $s \in S$ ,  $n \in \mathbb{N} + \{\infty\}$  and any policy  $\pi$  and it is not difficult to establish that

$$(3.4) \quad V_\pi(s) \leq E_{\pi_s} X_{n+1} \leq E_{\pi_s} X_n \leq v(s) \quad n \in \mathbb{N}.$$

Now let us consider a fixed  $\pi$  and  $s$  such that  $V_\pi(s) > -\infty$ . Then  $(X_n, n \in \mathbb{N} + \{\infty\})$  is a  $P_{\pi_s}$ -supermartingale relative to  $(\mathcal{F}_n, n \in \mathbb{N} + \{\infty\})$ . The limit

$$X'_\infty = \lim_n X_n$$

exists  $P_{\pi_s}$ -a.s.,  $(X_n, n \in \mathbb{N}, X'_\infty)$  is a  $P_{\pi_s}$ -supermartingale and  $X'_\infty \geq X_\infty$   $P_{\pi_s}$ -a.s. (cf. Neveu (1965), Proposition IV.5.4). Therefore we get

$$(3.5) \quad V_\pi(s) \leq E_{\pi_s} X'_\infty \leq E_{\pi_s} X_n \leq v(s) \quad n \in \mathbb{N}.$$

Perhaps the reader should be reminded that  $X'_\infty$  need not coincide with  $X_\infty$  (cf. also Theorem 6.3).

*Example 3.6.* Let  $S = \mathbb{N}$ ;  $A = \{0, 1\}$ ;  $D = S \times A$ ;  $q(1, \cdot, 1) = 1$ ,  $q(s, 0, s+1) = 1$  and  $q(s, 1, 1) = 1$  for all  $s > 1$ ;  $r(s, 0) = 0$  and  $r(s, 1) = 1 - 1/s$  for all  $s \in S$ . Then  $v := 1 \in \mathcal{E}_+$ . If  $f(s) = 0$ ,  $s \in S$  then for all  $s \geq 2$ ,  $X'_\infty = 1$  and  $X_\infty = 0$   $P_{f^\infty}$ -a.s.

For  $\tau \in T$  define  $X_\tau(h) = X_{\tau(h)}(h)$ ,  $h \in H$ .  $X_\tau$  is a random variable and  $E_{\pi_s} X_\tau^- \leq E_{\pi_s} X_\infty < \infty$  if  $V_\pi(s) > -\infty$ . The following theorem may be regarded as an *optional sampling theorem* for decision processes.

*Theorem 3.7* (cf. Hordijk (1974), Theorem 2.20). Let  $\pi \in \Delta$  and  $s \in S$  such that  $V_\pi(s) > -\infty$ . If  $\tau_1, \tau_2 \in T$  are stopping times with  $\tau_1 \leq \tau_2$ , then

$$E_{\pi_s} X_{\tau_1} \geq E_{\pi_s} X_{\tau_2}.$$

In particular,

$$v(s) \geq E_{\pi_s} X_\tau \geq V_\pi(s), \quad \tau \in T.$$

The *proof* follows from Neveu ((1965), Proposition IV.5.5), since the  $P_{\pi_s}$ -supermartingale  $(X_n, n \in \mathbb{N} + \{\infty\})$  is closed at the right.

#### 4. Optimal policies

Throughout this and the next two sections let  $\pi \in \Delta$  and  $s \in S$  be fixed such that  $-\infty < V_\pi(s) \leq v(s) < \infty$ . In the present paper we are concerned with the following concept of optimality.

Let  $\varepsilon \geq 0$ . The policy  $\pi$  is called  $(s, \varepsilon)$ -*optimal* (relative to  $v$ ) iff  $V_\pi(s) \geq v(s) - \varepsilon$ . The policy  $\pi$  is called  $(s, \varepsilon)$ -*thrifty* iff  $\lim_n E_{\pi s} X_n \geq v(s) - \varepsilon$ . The policy  $\pi$  is called  $(s, \varepsilon)$ -*equalizing* iff  $V_\pi(s) \geq \lim_n E_{\pi s} X_n - \varepsilon$ .  $\pi$  is called *s-optimal* (*s-thrifty*, *s-equalizing*) iff  $\pi$  is  $(s, 0)$ -optimal,  $((s, 0)$ -thrifty,  $(s, 0)$ -equalizing).

The notions ‘thrifty’ and ‘equalizing’ are adapted from Dubins and Savage (1965) (cf. also Sudderth (1972)). In this section we derive necessary and sufficient conditions for the policy  $\pi$  to be *s-optimal*. *s-thrifty* and *s-equalizing* policies are characterized in later sections.

As a consequence of (3.4) we obtain the following theorem.

**Theorem 4.1** (cf. Sudderth (1972), Theorem 8). Let  $\varepsilon \geq 0$ .  $\pi$  is  $(s, \varepsilon)$ -optimal iff  $\pi$  is  $(s, \varepsilon_1)$ -thrifty and  $(s, \varepsilon_2)$ -equalizing for some  $\varepsilon_1, \varepsilon_2$  such that  $\varepsilon_1 + \varepsilon_2 \leq \varepsilon$ .

**Theorem 4.2.** The following statements are equivalent:

- (i)  $\pi$  is *s-optimal*;
- (ii)  $\pi$  is *s-thrifty* and *s-equalizing*;
- (iii)  $(X_n, n \in \mathbf{N} + \{\infty\})$  is a  $P_{\pi s}$ -martingale;
- (iv)  $v(s) = E_{\pi s} X_\tau, \tau \in T$ .

*Proof.* (i)  $\Rightarrow$  (ii): obvious from Theorem 4.1. (ii)  $\Rightarrow$  (iii): this implication is easily deduced from (3.4) and the fact that  $(X_n, n \in \mathbf{N} + \{\infty\})$  is a  $P_{\pi s}$ -supermartingale. (iii)  $\Rightarrow$  (iv): Since the  $P_{\pi s}$ -martingale  $(X_n, n \in \mathbf{N} + \{\infty\})$  is closed at the right, the statement follows from Neveu ((1965), Proposition IV.5.5). (iv)  $\Rightarrow$  (i): immediate with  $\tau = \infty$ .

The following result may be regarded as a slight generalization of Theorem 4.2.

**Theorem 4.3.**  $\pi$  is *s-optimal* iff  $\pi$  is *s-thrifty* and  $\lim_n E_{\pi s} v^+ \circ \zeta_n = 0$ .

*Proof.* Since  $-\infty < V_\pi(s) \leq v(s) < \infty$  we get

$$(4.1) \quad \lim_n E_{\pi s} X_n = V_\pi(s) + \lim_n E_{\pi s} v^+ \circ \zeta_n.$$

In view of (4.1) the ‘if’ direction is obvious. The ‘only if’ direction is proved as follows. By Theorem 4.2 the policy  $\pi$  is *s-thrifty*. Now suppose that Assumption  $A^+$  holds. From Theorem 4.2 (iii) one obtains

$$E_{\pi s} v^+ \circ \zeta_n \leq E_{\pi s} \left[ \sum_{k=n}^{\infty} r_k^+ \right].$$

The passage to limit  $n \rightarrow \infty$  yields the desired assertion. If Assumption  $A^-$  is satisfied, the statement is derived from (4.1) and Remark 3.4.

Other criteria of optimality are based on the *optimality equation* (or Bellman

equation) and on the function  $V_\pi$  instead of  $v$ . These characterizations are contained in the next theorem. Note that the policy  $\pi$  is an arbitrary, not necessarily Markov, policy.

**Theorem 4.4.** Let  $v$  be the smallest excessive function in  $\mathcal{E}_i$ . Let  $\pi \in \Delta$  such that  $V_\pi(s) < \infty$ ,  $s \in S$ . The following statements are equivalent:

- (i)  $\pi$  is  $s$ -optimal for every  $s \in S$ ;
- (ii)  $V_\pi \in \mathcal{E}_i$ ;
- (iii)  $V_\pi \geq UV_\pi$  and  $\lim_n E_{\pi'} V_\pi \circ \zeta_n \geq 0$  on  $\{s: V_\pi(s) > -\infty\}$  for all  $\pi' \in \Delta$ .

*Proof.* The equivalence '(i)  $\Leftrightarrow$  (ii)' is obvious since  $v$  is the smallest function in  $\mathcal{E}_i$ . '(ii)  $\Leftrightarrow$  (iii)' follows from Remark 3.5.

In the *positive case*, i.e. if  $r \geq 0$ , the condition ' $\lim_n E_{\pi'} V_\pi \circ \zeta_n \geq 0$ ' is always satisfied.

## 5. Thrifty policies

We write

$$\varepsilon_1(s) = v(s) - E_{\pi s}[r_1 + v \circ \zeta_1]$$

and if  $n \geq 1$

$$\varepsilon_n(h) = v(s_n) - E_{\pi s}[r_n + v \circ \zeta_{n+1} | \chi_n = h], \quad h \in H_n.$$

Since  $v \in \mathcal{E}_i$  the function  $\varepsilon_n$  is  $(P_{\pi s})_{\chi_n}$ -a.s. non-negative, where  $(P_{\pi s})_{\chi_n}$  denotes the distribution of  $\chi_n$  under  $P_{\pi s}$ . Theorem 5.1 is a close cousin to Theorems 9 and 10 in Sudderth (1972).

**Theorem 5.1.** (a) Let  $\varepsilon \geq 0$ .  $\pi$  is  $(s, \varepsilon)$ -thrifty iff  $E_{\pi s}[\sum_1^\infty \varepsilon_n \circ \chi_n] \leq \varepsilon$ . (b) The following four statements are equivalent.

- (i)  $\pi$  is  $s$ -thrifty;
- (ii)  $v(s) = E_{\pi s} X_n$ ,  $n \in \mathbb{N}$ ;
- (iii)  $\varepsilon_n = 0$   $(P_{\pi s})_{\chi_n}$ -a.s.,  $n \in \mathbb{N}$ ;
- (iv)  $(X_n, n \in \mathbb{N})$  is a  $P_{\pi s}$ -martingale.

*Proof.* By induction on  $n$  we conclude  $E_{\pi s} X_n = v(s) - E_{\pi s} [\sum_{k=1}^n \varepsilon_k \circ \chi_k]$ . Then the Assertion (a) is obtained from the definition of  $(s, \varepsilon)$ -thrifty. The proof of (b) follows easily by (3.4) and Part (a).

Stated in terms of Markov or stationary policies, the theorem becomes as follows.

**Corollary 5.2.** Let  $\pi = (f_n)$  be a Markov policy.

- (a)  $\pi$  is thrifty iff  $v = L_{f_n} v$   $(P_{\pi s})_{\zeta_n}$ -a.s.,  $n \in \mathbb{N}$ .



(b) Suppose  $v = Uv$ . Then  $\pi$  is  $s$ -thrifty iff for any  $n \in \mathbf{N}$  and for  $(P_{\pi s})_{\zeta_n}$ -almost all  $s' \in S$  the point  $f_n(s')$  is a maximum point of the map

$$a \rightarrow [r(s', a) + \int q(s', a, dt)v(t)], \quad a \in D(s').$$

The criterion given in Corollary 5.2(b) is frequently used in dynamic programming problems. The additional assumption ' $v = Uv$ ' is required only for the 'if' direction.

**Corollary 5.3.** Let  $f \in F$ . Then  $f^\infty$  is  $s$ -thrifty for every  $s \in S$  if and only if  $v = L_f v$ .

Now we intend to generalize Theorem 5.1 by using stopping times. The following subset of  $T$  will be of interest. Let  $T_{\pi s}$  denote the set of all stopping times  $\tau \in T$  such that

$$(5.1) \quad \limsup_n \int_{\{\tau > n\}} X_n dP_{\pi s} \leq \int_{\{\tau = \infty\}} X_\infty dP_{\pi s}.$$

**Remark 5.4.** Since  $(X_n, n \in \mathbf{N} + \{\infty\})$  is a  $P_{\pi s}$ -supermartingale, (5.1) is equivalent to

$$\lim_n \int_{\{\tau > n\}} X_n dP_{\pi s} = \int_{\{\tau = \infty\}} X_\infty dP_{\pi s}.$$

Let us discuss some *properties* of the set  $T_{\pi s}$ .

Any bounded stopping time belongs to  $T_{\pi s}$ . If  $\tau \in T$  and  $P_{\pi s}(\tau < \infty) = 1$ , then  $\tau \in T_{\pi s}$  if and only if  $\lim_n \int_{\{\tau > n\}} X_n dP_{\pi s} = 0$ . In the *negative case*, i.e. if  $r \leq 0$ , the stopping time  $\tau = \infty$  belongs to  $T_{\pi s}$ , if  $v \leq 0$ . In the *absorbing case*, i.e. if there exists a measurable subset  $G$  of  $S$  such that  $q(s, a, S - G) = 0$ ,  $r(s, a) = 0$  for  $s \in G$ ,  $a \in D(s)$ , the stopping time  $\tau = \inf\{n \in \mathbf{N}: \zeta_n \in S - G\}$  belongs to  $T_{\pi s}$ , if  $v(s) = 0$  on  $G$ .

Let  $f \in F$  and  $q_f(s, s') = q(s, f(s), s')$ ,  $s, s' \in S$ . If there exists an absorbing set  $G_f$  (not necessarily independent of  $f$ ) relative to  $q_f$ , if  $r(s, a) \leq 0$  for  $s \in G_f$ ,  $a \in D(s)$  and if  $v(s) \leq 0$  for  $s \in G_f$ , then  $\tau = \inf\{n \in \mathbf{N}: \zeta_n \in S - G_f\}$  belongs to  $T_{\pi s}$ . For example, if  $S$  is countable  $G_f$  can be taken as the set of recurrent states of the Markov chain with transition law  $q_f$ .

We are now ready to state the main result of this section.

**Theorem 5.5.**  $\pi$  is  $s$ -thrifty iff  $v(s) = E_{\pi s} X_\tau$ ,  $\tau \in T_{\pi s}$ .

*Proof.* The 'if' direction follows from Theorem 5.1(b) with  $\tau = n \in T_{\pi s}$ . Let  $\pi$  be  $s$ -thrifty. By Theorem 5.1  $(X_n, n \in \mathbf{N})$  is a  $P_{\pi s}$ -martingale. Hence we infer

$$v(s) = \int_{\{\tau \leq n\}} X_\tau dP_{\pi s} + \int_{\{\tau > n\}} X_n dP_{\pi s} \quad n \in \mathbf{N}, \tau \in T.$$

Letting  $n \rightarrow \infty$  we get  $v(s) = E_{\pi s} X_\tau$  if  $\tau \in T_{\pi s}$ .

An immediate consequence of Theorem 5.5 is

**Theorem 5.6.** Let  $\pi$  be  $s$ -thrifty.  $\pi$  is  $s$ -optimal iff there exists a sequence  $(\tau_n)$  of stopping times  $\tau_n \in T_{\pi_s}$  such that  $\lim_n E_{\pi_s} X_{\tau_n} = V_{\pi}(s)$ .

We close this section with a useful characterization of the set  $T_{\pi_s}$ . As usual, we write  $\tau \wedge n = \min(\tau, n)$ .

**Theorem 5.7.** Let  $\pi$  be  $s$ -thrifty and  $\tau \in T$ . Then  $\tau \in T_{\pi_s}$  iff the sequence  $(X_{\tau \wedge n}^*)$  is uniformly  $P_{\pi_s}$ -integrable and  $X_s' = X_s$   $P_{\pi_s}$ -a.s. on  $\{\tau = \infty\}$ .

*Proof.* Since  $(X_n, n \in \mathbb{N}, X_s')$  is a  $P_{\pi_s}$ -supermartingale and  $X_s' \geq X_s$   $P_{\pi_s}$ -a.s. we get  $X_s' = X_s$  on  $\{\tau = \infty\}$ .  $\tau \wedge n$  belongs to  $T_{\pi_s}$  and thus by Theorem 5.5 we have

$$v(s) = E_{\pi_s} X_{\tau \wedge n} = E_{\pi_s} X_{\tau} \quad n \in \mathbb{N}.$$

Hence  $(X_{\tau \wedge n}, n \in \mathbb{N}, X_{\tau})$  is a  $P_{\pi_s}$ -martingale. By Neveu (1965), Proposition IV.5.1,  $(X_{\tau \wedge n})$  is uniformly  $P_{\pi_s}$ -integrable, and thus by definition  $(X_{\tau \wedge n}^*)$  is uniformly  $P_{\pi_s}$ -integrable. To show that the conditions are sufficient we apply Fatou's Lemma to the sequence  $(X_{\tau \wedge n})$  and obtain

$$\limsup_n E_{\pi_s} X_{\tau \wedge n} \leq E_{\pi_s} \left( \limsup_n X_{\tau \wedge n} \right) = E_{\pi_s} X_{\tau}$$

which completes the proof.

From the proof of Theorem 5.7 it follows that  $\tau \in T_{\pi_s}$  iff the sequence  $(X_{\tau \wedge n})$  is uniformly  $P_{\pi_s}$ -integrable and  $X_s' = X_s$   $P_{\pi_s}$ -a.s. on  $\{\tau = \infty\}$ .

## 6. Equalizing policies

In negative dynamic programming problems every policy is  $s$ -equalizing. More generally, the condition (EN) of Hinderer (1970), the condition  $(C^+)$  of Hinderer (1971) and the condition (C) of Schäl (1975) guarantee that all policies are equalizing. (This can be proved by Theorem 6.1.) In positive problems, this is not necessarily true (not even for thrifty policies), e.g. in Example 3.6 the stationary policy  $f^*$  is thrifty, but not  $s$ -equalizing for  $s > 1$ . However, if  $S$  is finite and if Assumption  $A^+$  holds then every  $s$ -thrifty stationary policy is  $s$ -equalizing.

**Theorem 6.1.** (a) Let  $\varepsilon \geq 0$ .  $\pi$  is  $(s, \varepsilon)$ -equalizing iff  $\lim_n E_{\pi_s} v \circ \zeta_n \leq \varepsilon$ . (b)  $\pi$  is  $s$ -equalizing iff  $\tau = \infty$  belongs to  $T_{\pi_s}$ .

*Proof.* By Remark 3.3 (a) follows from (4.1). (b) is obvious by Remark 5.4.

Theorem 6.2 generalizes Theorem 4.10 in Chow, Robbins and Siegmund (1971) to more general dynamic programming problems.

**Theorem 6.2.** Let  $\pi$  be  $s$ -thrifty. Then  $\pi$  is  $s$ -optimal iff the sequence  $(X_n^*)$  is uniformly  $P_{\pi s}$ -integrable and  $X_\infty' = X_\infty$   $P_{\pi s}$ -a.s.

The *proof* is a consequence of Theorems 4.2, 5.7 and 6.1(b).

If the policy  $\pi$  is  $s$ -equalizing then  $V_\pi(s) = E_{\pi s} X_\infty = E_{\pi s} X_\infty'$  by (3.5). The converse fails in general. But the converse is true if the sequence  $(v^+ \circ \zeta_n)$  is uniformly  $P_{\pi s}$ -integrable. This additional condition is satisfied in Sudderth (1972). Note that  $\lim_n v \circ \zeta_n \geq 0$   $P_{\pi s}$ -a.s. exists.

**Theorem 6.3.** Let the sequence  $(v^+ \circ \zeta_n)$  be uniformly  $P_{\pi s}$ -integrable. Then the following statements are equivalent:

- (i)  $\pi$  is  $s$ -equalizing;
- (ii)  $V_\pi(s) = E_{\pi s} X_\infty'$ ;
- (iii)  $P_{\pi s}(X_\infty = X_\infty') = 1$ ;
- (iv)  $\lim_n v \circ \zeta_n = 0$   $P_{\pi s}$ -a.s.

*Proof.* From (3.5) and Fatou's Lemma we conclude

$$(6.1) \quad \lim_n E_{\pi s} v \circ \zeta_n = E_{\pi s} \left( \lim_n v \circ \zeta_n \right).$$

Then the statements are immediate.

We remark that Assumption  $A^+$  need not imply that the sequence  $(v^+ \circ \zeta_n)$  is uniformly  $P_{\pi s}$ -integrable. The following theorem yields a slight modification of Theorems 11, 12, and 13 in Sudderth (1972).

**Theorem 6.4.** Let the sequence  $(v^+ \circ \zeta_n)$  be uniformly  $P_{\pi s}$ -integrable and let  $\varepsilon > 0$ .

- (a)  $\pi$  is  $(s, \varepsilon)$ -equalizing iff  $E_{\pi s} X_\infty' \leq V_\pi(s) + \varepsilon$ .
- (b) If  $P_{\pi s}(\lim_n v \circ \zeta_n \leq \varepsilon) = 1$  then  $\pi$  is  $(s, \varepsilon)$ -equalizing.
- (c) If  $\pi$  is  $(s, \varepsilon^2)$ -equalizing then  $P_{\pi s}(\lim_n v \circ \zeta_n \leq \varepsilon) \geq 1 - \varepsilon$ .

*Proof.* From (6.1) we get  $\lim_n E_{\pi s} X_n = E_{\pi s} X_\infty'$ . Then (a) follows by definition. In view of (a) the proof of (b) (resp. (c)) is similar to the proof of Theorem 12 (resp. Theorem 13) of Sudderth (1972) and is therefore omitted.

## 7. Applications

Throughout this section,  $S$  is assumed to be countable. Then the value-function  $V$  belongs to  $\mathcal{E}_i$ . Here we suppose  $v = V$ . A policy  $\pi$  is called *optimal* if  $V_\pi(s) = V(s)$  for every  $s \in S$ .

The following result has been proved by Blackwell (1970), Ornstein (1969) and Hordijk (1974) in the positive case.

**Theorem 7.1.** Assume  $A^+$ . If there exists an optimal Markov policy then there exists an optimal stationary policy.

*Proof.* Let  $\pi = (f_n)$  be the optimal Markov policy. Without loss of generality we can assume that  $V_\pi(s) = V(s) > -\infty$ ,  $s \in S$ . Let us consider the positive stationary dynamic program  $(S, (A, \mathcal{A}), D', q, r')$  with

$$D' = \{(s, a) \in D : V(s) = r(s, a) + \sum_{t \in S} q(s, a, t) V(t)\}$$

$$r'(s, a) = r^+(s, a), (s, a) \in D'.$$

In virtue of Theorem 4.2 and Corollary 5.2 the graph of  $f_1$  belongs to  $D'$  and thus the dynamic program above is well-defined. Without loss of generality we can assume in view of Corollary 5.2 that  $\pi$  belongs to  $\Delta'$ . Therefore we conclude

$$(7.1) \quad V(s) \leq V'(s) \leq W_+(s) < \infty \quad s \in S.$$

Let  $\varepsilon > 0$ . According to a theorem of Ornstein (1969) there exists a stationary policy  $f^\varepsilon$  such that

$$(7.2) \quad f(s) \in D'(s), \quad s \in S$$

$$(7.3) \quad V_{f^\varepsilon}(s) \geq (1 - \varepsilon) V'(s), \quad s \in S.$$

Corollary 5.3 and (7.2) tell us that  $f^\varepsilon$  is  $s$ -thrifty for  $s \in S$ . By (7.1) and (7.3) we get  $E_{f^\varepsilon} V \circ \zeta_n \leq (1 - \varepsilon)^{-1} E_{f^\varepsilon} V_{f^\varepsilon} \circ \zeta_n \rightarrow 0 (n \rightarrow \infty)$  and by Theorem 6.1,  $f^\varepsilon$  is  $s$ -equalizing for  $s \in S$ . Consequently  $f^\varepsilon$  is optimal.

Now we consider a stationary policy  $f^*$ . Theorem 5.6 enables us to state sufficient conditions for the  $s$ -optimality of  $f^*$ .

**Theorem 7.2.** Assume  $A^+$ . Let  $f^*$  be  $s$ -thrifty and  $V_{f^*}(s) > -\infty$ . If there exists a sequence  $(\tau_n)$  of stopping times  $\tau_n \in T_{f^*}$  such that  $\tau_n \geq n$  and

$$(7.4) \quad \limsup_n E_{f^*} [V \circ \zeta_{\tau_n} 1_{N \circ \tau_n}] \leq 0$$

then  $f^*$  is  $s$ -optimal, i.e.  $V_{f^*}(s) = V(s)$ .

*Proof.* Assumption  $A^+$  implies  $V(s) < \infty$ . From (7.4) and Fatou's Lemma one obtains  $\lim_n E_{f^*} X_{\tau_n} \leq V_{f^*}(s)$ . The assertion follows from Theorems 5.5 and 5.6.

The Condition (7.4) is satisfied e.g. if there exists an absorbing subset  $G_f$  of  $S$  such that  $r(s, a) \leq 0$  for  $s \in G_f$ ,  $a \in D(s)$

$$\sup_{j \notin G_f} V(j) < \infty \quad \text{and} \quad \lim_n P_{f^*}(\tau_n < \infty) = 0,$$

where  $\tau_n = \inf\{k \geq n : \zeta_k \in S - G_f\}$  (cf. Section 5).

## References

- BLACKWELL, D. (1965) Discounted dynamic programming. *Ann. Math. Statist.* **36**, 226–235.
- BLACKWELL, D. (1970) On stationary policies. *J. R. Statist. Soc. A* **133**, 33–37.
- CHOW, Y. S., ROBBINS, H. AND SIEGMUND, D. (1971) *Great Expectations: The Theory of Optimal Stopping*. Houghton-Mifflin Company, Boston.
- DUBINS, L. E. AND SAVAGE, L. J. (1965) *How to Gamble if You Must: Inequalities for Stochastic Processes*. McGraw-Hill, New York.
- HINDERER, K. (1970) Foundations of non-stationary dynamic programming with discrete time parameter. *Lecture Notes in Operations Research and Mathematical Systems*. **33**, Springer, Berlin.
- HINDERER, K. (1971) Instationäre dynamische Optimierung bei schwachen Voraussetzungen über die Gewinnfunktionen. *Abh. Math. Sem. Univ. Hamburg* **36**, 208–223.
- HORDIJK, A. (1974) *Dynamic Programming and Markov Potential Theory*. Mathematical Centre Tracts No. 51, Amsterdam.
- NEVEU, J. (1965) *Mathematical Foundations of the Calculus of Probability*. Holden-Day, San Francisco.
- ORNSTEIN, D. (1969) On the existence of stationary optimal strategies. *Proc. Amer. Math. Soc.* **20**, 563–569.
- RIEDER, U. (1975) On stopped decision processes with discrete time parameter. *Stoch. Proc. Appl.* **3**, 365–383.
- ROSS, S. M. (1974) Dynamic programming and gambling models. *Adv. Appl. Prob.* **6**, 593–600.
- SCHÄL, M. (1975) Conditions for optimality in dynamic programming and for the limit of  $n$ -stage optimal policies to be optimal. *Z. Wahrscheinlichkeitsth.* **32**, 179–196.
- STRAUCH, R. E. (1966) Negative dynamic programming. *Ann. Math. Statist.* **37**, 871–890.
- SUDDERTH, W. (1972) On the Dubins and Savage characterization of optimal strategies. *Ann. Math. Statist.* **43**, 498–507.