

# A NOTE ON UNDISCOUNTED DYNAMIC PROGRAMMING<sup>1</sup>

BY ASHOK MAITRA<sup>2</sup>

*Mathematisch Centrum, Amsterdam*

**1. Introduction.** We consider a system with a finite number of states  $1, 2, \dots, S$ . Once a day, we observe the current state  $s$  of the system and choose an action  $a$  from an arbitrary set  $A$  of actions. As a result, two things happen: (1) we receive an immediate income  $i(s, a)$ , and (2) the system moves to a new state  $s'$  with probability  $q(s' | s, a)$ . Assume that the incomes are bounded, that is, there exists a positive number  $M$  such that  $|i(s, a)| \leq M$ ,  $s = 1, 2, \dots, S$ ,  $a \in A$ . The problem is to maximise the average rate of income (to be defined below).

Denote by  $F$  the set of all functions  $f$  on  $S$  into  $A$ . A *policy*  $\pi = \{f_1, f_2, \dots\}$  is a sequence of functions  $f_n \in F$ . Thus, to use policy  $\pi$  is to choose the action  $f_n(s)$  on the  $n$ th day, if the system is in state  $s$  on that day. We shall call a policy  $\pi = \{f_n\}$  *stationary* if  $f_n = f$ ,  $n = 1, 2, \dots$ , and denote it by  $f^{(\infty)}$ .

With each  $f \in F$ , associate (1) the  $S \times 1$  vector  $r(f)$ , whose  $s$ th coordinate is  $i(s, f(s))$  and (2) the  $S \times S$  stochastic matrix  $Q(f)$ , whose  $(s, s')$  element is  $q(s' | s, f(s))$ . Hence, if we use the policy  $\pi = \{f_n\}$ , the  $n$ -step transition matrix of the system is  $Q_n(\pi) = \prod_{k=1}^n Q(f_k)$ . In particular, if our policy is stationary, the system becomes a discrete time-parameter Markov chain with stationary transition probabilities.

Given a policy  $\pi$ , let us denote by  $W_n(\pi)$  the  $S \times 1$  vector of incomes on the  $n$ th day, when the policy  $\pi$  is used. Set

$$x(\pi) = \lim_{N \rightarrow \infty} N^{-1} \sum_{n=1}^N W_n(\pi)$$

whenever the limit exists. Blackwell [1] has shown that the limit exists whenever  $\pi$  is a stationary policy. In the case of a stationary policy,  $x(f^{(\infty)})$  is the vector of average rates of income, when the policy  $f^{(\infty)}$  is used.

We shall say that a policy  $f_0^{(\infty)}$  is *optimal* among stationary policies if  $x(f_0^{(\infty)}) \geq x(f^{(\infty)})$  for all  $f \in F$  (for any two  $S \times 1$  vectors  $w_1$  and  $w_2$ , we shall write  $w_1 \geq w_2$  if every coordinate of  $w_1$  is at least as large as the corresponding coordinate of  $w_2$ , and  $w_1 > w_2$  if  $w_1 \geq w_2$  and  $w_1 \neq w_2$ ).

Blackwell [1] showed that, if  $A$  is finite, there exists an optimal policy among stationary policies. When  $A$  is not finite, there may not exist an optimal policy. Consider, for instance, a system with a single state and  $A = \{1, 2, \dots\}$ . Choice of action  $i$  brings an income of  $1 - 1/i$  dollars. It is clear that there is no optimal stationary policy.

The purpose of this note is to prove:

**THEOREM.** *Let  $A$  be arbitrary. Given  $\epsilon > 0$ , there exists a stationary policy  $f_\epsilon^{(\infty)}$*

Received 24 November 1965.

<sup>1</sup> Report SP 89 of the Statistics Department, Mathematisch Centrum, Amsterdam.

<sup>2</sup> Now with Indian Statistical Institute, Calcutta.