

**U.S. DEPARTMENT OF COMMERCE  
National Technical Information Service**

**AD-A032 230**

**A Model for the Analysis of Markovian  
Decision Processes with Unobservable  
States and Unobservable Costs**

**Rand Corp Santa Monica Calif**

**Oct 72**

ADA032230

REPRODUCED BY  
NATIONAL TECHNICAL  
INFORMATION SERVICE  
U. S. DEPARTMENT OF COMMERCE  
SPRINGFIELD, VA. 22161

A MODEL FOR THE ANALYSIS OF MARKOVIAN DECISION PROCESSES  
WITH UNOBSERVABLE STATES AND UNOBSERVABLE COSTS

James D. Steele\*  
The Rand Corporation, Santa Monica, California

INTRODUCTION

Consider the Markovian Decision Process (MDP) defined by the following objects;

State Space  $S = \{1, 2, 3, \dots, N\}$  , for finite  $N$ ,

Action Space  $A = \{a_1, a_2, \dots, a_M\}$  , for finite  $M$ ,

Cost Set  $C = \{C(i, a_j) : i \in S, a_j \in A\}$  ,

Transition Probabilities =  $\{q_{ij}(a_k) : i, j \in S, a_k \in A\}$  ,

Discount Factor  $\alpha$ , such that  $0 < \alpha < 1$ .

The problem is to find a policy for taking actions which minimizes the total expected discounted cost over the infinite future, given the initial state of the process.

A stationary policy for (MDP) is defined as a map  $f : S \rightarrow A$ . Howard [2] analyzed (MDP's) having finite state and action spaces and proved that an optimal stationary policy (i.e. a stationary policy which minimizes the total expected discounted cost) always exists. The Howard Policy Improvement Routine is a method by which an optimal stationary policy for (MDP) may be found.

Suppose now, that we are given the (MDP) as defined above, but that we are not allowed to observe the state at any observation point  $t = 0, 1, 2, \dots$ . Suppose also, that we are not allowed to observe

---

\* The author is an Applied Mathematician at The Rand Corporation.

the cost  $C(X_t, a_t)$  at any observation points  $t = 0, 1, 2, \dots$ . In other words, the total cost will be assessed at infinity. Finally, suppose that we are allowed to observe the initial probability distribution over the state space  $S$ . In this paper we develop a model for analyzing this problem and present some preliminary results. A rather thorough treatment of the problem of unobservable states and unobservable costs for (MDP's) having two states and two actions may be found in Steele [3].

### THE MODEL

In an effort to analyze the above problem, we define the following objects.

$$\begin{aligned} \mathcal{S} &= \{ \text{All probability distributions over } S \} \\ &= \{ \bar{P} = (P_1, P_2, \dots, P_N) \in E_N : 0 \leq P_i \leq 1, \sum_{i=1}^N P_i = 1, \\ &\quad i=1, 2, \dots, N \}, \end{aligned}$$

where  $E_N$  is  $N$ -dimensional Euclidean space and we let  $P_i$  be the probability of being in state  $i$ ;

the set  $\mathcal{A} = A = \{a_1, a_2, \dots, a_M\}$ ; the transition

matrices :  $Q(a_K) = [q_{ij}(a_K)]$ ; and the cost vectors

$$\bar{C}(a_K) = (C[1, a_K], \dots, C[N, a_K]).$$

We note that if the distribution over  $S$  is  $\bar{P} \in \mathcal{S}$  and we take action  $a_K \in \mathcal{A}$ , then the new distribution over  $S$  will be given by  $\bar{P}' = \bar{P}Q(a_K)$ . The expected cost,  $\mathcal{C}(\bar{P}, a_K)$ , of having the distribution  $\bar{P}$  and taking action  $a_K$  will be given as the inner product

$$\mathcal{C}(\bar{P}, a_K) = \langle \bar{P}, \bar{C}(a_K) \rangle = \sum_{i=1}^N P_i C(i, a_K).$$

At this point, we note that the new distribution  $\bar{P}'$  depends only on the current distribution  $\bar{P}$  and the current action  $a_K$ , i.e.  $\bar{P}' = \bar{P}Q(a_K)$ , therefore, we see that we now have a new Markov an Decision Process (MDP) defined by the objects

State Space  $\tilde{S} = \{\text{all probability distributions over } S\}$ ,

Action Space  $\tilde{A} = A = \{a_1, a_2, \dots, a_M\}$ , and

Cost Set  $\tilde{C} = \{C(\bar{P}, a_K) : \bar{P} \in \tilde{S}, a_K \in \tilde{A}\}$

Discount Factor  $\alpha$ , such that  $0 < \alpha < 1$ .

The set of all stationary policies for (MDP) is given by  $F = \{f : \tilde{S} \rightarrow \tilde{A} = A\}$ . For any such  $f \in F$  and initial state  $\bar{P}_0 \in \tilde{S}$ , the total expected discounted cost is given by

$$\begin{aligned} V_f(\bar{P}_0) &= \sum_{t=0}^{\infty} \alpha^t C(\bar{P}_t, f[\bar{P}_t]) \\ &= \sum_{t=0}^{\infty} \alpha^t \langle \bar{P}_t, \bar{C}[f(\bar{P}_t)] \rangle, \text{ where} \end{aligned}$$

$$\bar{P}_t = \bar{P}_0 Q(f[\bar{P}_0]) Q(f[\bar{P}_1]) \dots Q(f[\bar{P}_{t-1}]) \text{ for } t=1, 2, 3, \dots$$

The (MDP) as defined above (having uncountable state space and finite action space) belongs to the class of problems analyzed by Blackwell [1]. His analysis showed that an optimal stationary policy always exists and that the Howard Policy Improvement Routine may be extended to this problem. However, in the finite state-finite action problems the set of all stationary policies is finite and, therefore, the Howard Policy Improvement Routine will produce an optimal stationary policy in a finite number of steps. In the uncountable state-finite action problem, the set of all stationary policies is uncountable and, therefore, the Howard Policy Improvement Routine cannot, in general, be used as a method for actually finding an optimal policy.

CONVEX-STATIONARY POLICIES

With  $F$  as the set of all stationary policies for  $(M^0P)$ , we define the set  $\mathcal{F} \subset F$ , of all convex-stationary policies for  $(M^0P)$ , as

$$\mathcal{F} = \{f \in F : f^{-1}(a_k) \text{ is a convex set for each } a_k \in A\}.$$

CONSTANT SEQUENCE POLICIES

Given the action space  $A = \{a_1, a_2, \dots, a_M\}$ , we define  $A^N = A \times A \times \dots \times A$ ,  $N$ -factors,  $N = 1, 2, 3, \dots$ , to be the set of all sequences of length  $N$  of elements of  $A$ , and we define  $A^\infty = A \times A \times A \times \dots$ , to be the set of all infinite sequences of elements of  $A$ . For any finite sequence  $S \in A^N$ ,  $N \geq 1$ , we define the sequence  $S^K \in A^{NK}$  to be the sequence  $S^K = s, s, \dots, s$ ,  $K$ -factors, and  $S^\infty \in A^\infty$ , to be the sequence  $S^\infty = s, s, s, \dots$ . For any finite sequences  $S_1$  and  $S_2$  we define  $A(S_1; S_2) = A\{S_1, S_2\}$  to be the set consisting of the two action sequences  $S_1$  and  $S_2$  and

$A(S_1; S_2)^\infty = A(S_1; S_2) \times A(S_1; S_2) \times \dots$ , to be the set of all infinite sequences of elements of  $A(S_1; S_2)$ . For any finite sequence  $S = a_1, a_2, \dots, a_N$ ,  $N \geq 1$ ,  $a_i \in A$ , we define

$$L_S(\bar{P}) = \sum_{t=0}^{N-1} \alpha^t C(\bar{P}_t, a_{t+1}),$$

for  $\bar{P} \in \mathcal{S}$  and  $\bar{P}_0 = \bar{P}$ , to be the cost of starting at  $\bar{P}$  and operating for  $N$  time periods when using the  $i$ -th entry in  $S$ ,  $1 \leq i \leq N$ , as the action to be taken at the  $i$ -th observation time. We say that  $L_S(\cdot)$  is the cost of using the finite sequence  $S$ , for all initial  $\bar{P} \in \mathcal{S}$ . If  $S \in A^\infty$ , we define the constant sequence policy  $(S)$ , to be the policy which uses the sequence  $S$  when starting at any initial  $\bar{P} \in \mathcal{S}$ , we define  $(A^\infty)$  to be the set of all such policies and we use  $V_{(S)}(\cdot)$ , in place of  $L_{(S)}(\cdot)$ , for the cost of the policy  $(S)$ .

IMMEDIATE RESULTS

LEMMA 1: The cost function,  $V_{(r)}$ , for any policy  $(r) \in (A^\infty)$  is linear on  $\mathcal{Z}$ .

Proof: Let the sequence  $r \in A^\infty$  be given by  $r = a_1, a_2, a_3, \dots$ . Now, for any points  $\bar{P}, \bar{P}', \bar{P}''$  in  $\mathcal{Z}$  such that  $\bar{P} = \lambda \bar{P}' + (1-\lambda) \bar{P}''$  for some  $\lambda \in [0, 1]$  we have (with  $Q(a_0) \equiv$  the Identity operator),

$$\begin{aligned} V_{(r)}(\bar{P}) &= \sum_{t=0}^{\infty} \alpha^t \bar{c}(\bar{P}_t, a_{t+1}) \\ &= \sum_{t=0}^{\infty} \alpha^t \langle \bar{P}Q(a_0) Q(a_1) \dots Q(a_t), \bar{c}(a_{t+1}) \rangle \\ &= \lambda \sum_{t=0}^{\infty} \alpha^t \langle \bar{P}'Q(a_0) \dots Q(a_t), \bar{c}(a_{t+1}) \rangle \\ &\quad + (1-\lambda) \sum_{t=0}^{\infty} \alpha^t \langle \bar{P}''Q(a_0) \dots Q(a_t), \bar{c}(a_{t+1}) \rangle \\ &= \lambda V_{(r)}(\bar{P}') + (1-\lambda) V_{(r)}(\bar{P}'') \end{aligned}$$

or

$$V_{(r)}(\lambda \bar{P}' + [1-\lambda] \bar{P}'') = \lambda V_{(r)}(\bar{P}') + (1-\lambda) V_{(r)}(\bar{P}''). \quad \text{Q.E.D.}$$

For any stationary policy  $f$  and any  $\bar{P} \in \mathcal{Z}$ , we define the sequence  $S(\bar{P}, f) \in A^\infty$ ,

$$S(\bar{P}, f) = \{a_t\}, \text{ by } a_{t+1} = f(\bar{P}_t),$$

$$\bar{P}_{t+1} = \bar{P}_t Q(f[\bar{P}_t]), \text{ for } t = 0, 1, 2, \dots, \text{ and } \bar{P}_0 = \bar{P}.$$

We say that  $S(\bar{P}, f)$  is the sequence generated at  $\bar{P}$  when the stationary policy  $f$  is used.

LEMMA 2: For any stationary policy  $f$  and any  $\bar{P} \in \mathcal{Z}$ , we have

$J_f(\bar{P}) = V_{(S[\bar{P}, f])}(\bar{P})$ , where  $V_f(\bar{P})$  is the cost of using the stationary policy  $f$  and starting at  $\bar{P}$ .

Proof: By definition of  $S(\bar{P}, f)$ . Q.E.D.

THEOREM 1: The optimal cost function is concave on  $\mathcal{S}$ .

Proof: Let  $f^*$  be optimal. Lemmas 1 and 2 show that at each  $\bar{P} \in \mathcal{S}$

$$V_{f^*}(\bar{P}) = V_{(S[\bar{P}, f^*])}(\bar{P}) = \inf_{S \in A} V(S)(\bar{P})$$

Therefore, we see that at each point  $\bar{P} \in \mathcal{S}$ ,  $V_{f^*}(\bar{P})$  is the infimum over a set of linear functions and hence  $V_{f^*}$  is concave. Q.E.D.

Next we prove that the optimal cost function is continuous on  $\mathcal{S}$  by making use of the following representation. Let  $B$  be the set of all bounded Baire functions on  $\mathcal{S}$ . Define a norm on  $B$  by

$$\|V\| = \sup_{\bar{P} \in \mathcal{S}} |V(\bar{P})|, \text{ for any } V \in B.$$

Next, define the operator  $U : B \rightarrow B$  by

$$(Uv)(\bar{P}) = \min_{a_K \in A} \{L_{a_K}(\bar{P}) + \alpha V(PQ[a_K])\}.$$

In Lemma 3, we state some results presented in Reference [1].

LEMMA 3.

- (i)  $U$  is a contraction operator
- (ii) For any  $V \in B$ , the sequence  $U_n = U^n V$  converges to the optimal cost function  $V_{f^*}$ .
- (iii) The optimal cost function,  $V_{f^*}$ , is the unique solution to  $UV_{f^*} = V_{f^*}$ .

We now have the following Theorem.

THEOREM 2: The optimal cost function is uniformly continuous on  $\bar{S}$ .

Proof: For any  $u \in B$ , we have

$$u_n = U^n u \rightarrow V_{f^*} \text{ as } n \rightarrow \infty.$$

We also have

$$u_{n+1}(\bar{P}) = \min_{a_k \in A} \{L_{a_k}(\bar{P}) + \alpha u_n(\bar{P}Q[a_k])\}$$

for  $\bar{P} \in S$ . Therefore, we see that since  $L_{a_k}(\cdot)$  is continuous for each  $a_k \in A$ , each  $u_n$  will be continuous if  $u = u_0$  is continuous. The convergence of  $u_n \rightarrow V_{f^*}$  is uniform because  $U$  is a contraction operator, i.e.

$$\|u_n - V_{f^*}\| \leq \alpha^n \|u - V_{f^*}\|. \quad \text{Q.E.D.}$$

BIBLIOGRAPHY

- [1] Blackwell, D. (1965). "Discounted Dynamic Programming." Annals of Mathematical Statistics 36, pp. 226-235.
- [2] Howard, R.A. (1960). Dynamic Programming and Markov Processes. Cambridge, Mass.; Technology Press, and New York: Wiley.
- [3] Steele, J.D., "On Two-state Two-action Markovian Decision Processes With Unobservable States and Unobservable Costs." Unpublished Ph.D. dissertation, Department of Mathematics, University of California at Los Angeles, 1972.