

DOCUMENT RESUME

ED 124 161

IR 003 551

AUTHOR Wollmer, Richard D.
 TITLE Partially Observable Markov Decision Processes Over
 an Infinite Planning Horizon with Discounting.
 Technical Report No. 77.

INSTITUTION University of Southern California, Los Angeles.
 SPONS AGENCY Behavioral Technology Labs.
 Advanced Research Projects Agency (DOD), Washington,
 D.C.; Office of Naval Research, Washington, D.C.
 Personnel and Training Research Programs Office.

PUB DATE Mar 76
 CONTRACT N0014-75-C-0838
 NOTE 25p.

EDRS PRICE MP-\$0.83 HC-\$1.67 Plus Postage.
 DESCRIPTORS Computer Assisted Instruction; *Decision Making;
 Instructional Systems; *Linear Programming;
 Mathematical Applications; *Mathematical Models;
 Operations Research; Probability Theory; *Systems
 Approach
 IDENTIFIERS Markov Models

ABSTRACT

The true state of the system described here is characterized by a probability vector. At each stage of the system an action must be chosen from a finite set of actions. Each possible action yields an expected reward, transforms the system to a new state in accordance with a Markov transition matrix, and yields an observable outcome. The problem of finding the total maximum discounted reward as a function of the probability state vector may be formulated as a linear program with an infinite number of constraints. The reward function may be expressed as a partial N-dimensional Maclaurin series. The coefficients in this series are also determined as an optimal solution to a linear program with an infinite number of constraints. A sequence of related finitely constrained linear programs is solved which then generates a sequence of solutions that converge to a local minimum for the infinitely constrained program. This model is applicable to computer assisted instruction systems as well as to other situations. (Author/CH)

 * Documents acquired by ERIC include many informal unpublished *
 * materials not available from other sources. ERIC makes every effort *
 * to obtain the best copy available. Nevertheless, items of marginal *
 * reproducibility are often encountered and this affects the quality *
 * of the microfiche and hardcopy reproductions ERIC makes available *
 * via the ERIC Document Reproduction Service (EDRS). EDRS is not *
 * responsible for the quality of the original document. Reproductions *
 * supplied by EDRS are the best that can be made from the original. *

U.S. DEPARTMENT OF HEALTH,
EDUCATION & WELFARE
NATIONAL INSTITUTE OF
EDUCATION

THIS DOCUMENT HAS BEEN REPRO-
DUCED EXACTLY AS RECEIVED FROM
THE PERSON OR ORGANIZATION ORIGIN-
ATING IT. POINTS OF VIEW OR OPINIONS
STATED DO NOT NECESSARILY REPRE-
SENT OFFICIAL NATIONAL INSTITUTE OF
EDUCATION POSITION OR POLICY.

DEPARTMENT OF PSYCHOLOGY
UNIVERSITY OF SOUTHERN CALIFORNIA

ED124161

Technical Report No. 77

PARTIALLY OBSERVABLE MARKOV DECISION PROCESSES
OVER AN INFINITE PLANNING HORIZON
WITH DISCOUNTING

March 1976

Richard D. Wollmer

Sponsored by

Personnel and Training Research Programs
Psychological Sciences Division
Office of Naval Research

and

Advanced Research Project Agency
Under Contract No. N00014-75-C-0838

The views and conclusions contained in this document
are those of the authors and should not be interpreted
as necessarily representing the official policies,
either expressed or implied of the Office of Naval
Research, the Advanced Research Projects Agency, or
the U.S. Government.

Approved for public release: distribution unlimited.

55
R083
H

ARPA TECHNICAL REPORT

March 1976

1. ARPA Order Number	:	2284
2. ONR NR Number	:	154-355
3. Program Code Number	:	1 B 729
4. Name of Contractor	:	University of Southern California
5. Effective Date of Contract	:	76 January 1
6. Contract Expiration Date	:	76 December 31
7. Amount of Contract	:	\$200,000.00
8. Contract Number	:	N00014-75-C-0838
9. Principal Investigator	:	J. W. Rigney (213) 746-2127
10. Scientific Officer	:	Marshall Farr
11. Short Title	:	Learning Strategies

This Research Was Supported

by

The Advanced Research Projects Agency

and by

The Office of Naval Research

And Was Monitored by

The Office of Naval Research

SUMMARY

This is the last in a series of technical reports concerned with mathematical approaches to instructional sequence optimization in instructional systems. The problem treated here is very closely related to that treated by Smallwood and Sondik (4). Both papers deal with Markov decision processes where the true state of the system is not known with certainty. Hence the state of the system is characterized by a probability vector. Each action yields an expected reward, transforms the system to a new state and yields an observable outcome. One wishes to determine an action for each probability state vector so as to maximize the total expected reward. Smallwood and Sondik (4) solve this problem exactly for a finite time horizon. This report treats the infinite time horizon with a discount factor, using a partial N dimensional Maclaurin series to approximate the total optimal reward as a function of the probability state vector. While this model was developed for computed aided instruction, it is applicable to other situations as well. This model also is of considerable theoretical value.

ABSTRACT

This paper describes a system that may be in any one of states $1, 2, \dots, N$. The true state of the system is not known with certainty, and consequently is described by a probability vector. At each stage an action must be chosen from a finite set. Each possible action returns an expected reward, transforms the system to a new state in accordance with a Markov transition matrix, and yields an observable outcome. It is required to determine an action for each possible state vector in order to maximize the total expected reward over an infinite time horizon under a discount factor, β , where $0 < \beta < 1$.

The problem of finding the total maximum discounted reward as a function of the probability state vector may be formulated as a linear program with an infinite number of constraints. The reward function may be expressed as an N dimensional Maclaurin series and in this paper it is approximated by a partial series consisting of terms up to degree n . The coefficients in this series are also determined as an optimal solution to a linear program with an infinite number of constraints. A sequence of related finitely constrained linear programs are solved which generate a sequence of solutions that converge to a local minimum for the infinitely constrained program. It is an open question as to whether this local minimum is actually a global minimum. However it should be noted that the function being approximated is convex and consequently has the property that any local minimum is a global one as well.

ACKNOWLEDGEMENTS

The research discussed in this report was monitored by Dr. Marshall Farr and Dr. Joseph Young, Personnel and Training Research Programs, Office of Naval Research, and Dr. Harry F. O'Neil, Jr., Program Manager, Human Resources Research, Advanced Research Projects Agency. Their support and encouragement is gratefully acknowledged.

TABLE OF CONTENTS

<u>Section</u>	<u>Page</u>
I. INTRODUCTION	1
II. FORMULATION	2
III. A LEARNING EXAMPLE	4
IV. THE MAXIMUM REWARD FUNCTION	7
V. LINEAR PROGRAM FORMULATION	9
VI. COMPUTATIONAL PROCEDURE	13
VII. BOUNDS ON ACCURACY	17
REFERENCES	18

PARTIALLY OBSERVABLE MARKOV DECISION
PROCESSES OVER AN INFINITE PLANNING
HORIZON WITH DISCOUNTING

1. Introduction

This paper describes a system that may be in anyone of states $1, 2, \dots, N$. The true state of the system is not known with certainty and consequently is described by a probability vector. At each stage an action must be chosen from a finite set. This action returns an expected reward, transforms the system to a new (but not necessarily different) state according to a Markov process, and yields an observable outcome. The problem addressed here is that of determining an action for each possible state vector in order to maximize the total expected reward over an infinite horizon under a discount factor, β , where $0 < \beta < 1$.

Smallwood and Sondik (4) have treated this problem for the finite horizon case without a discount factor and have determined that the total maximum expected reward is a piecewise linear function of the probability state vector. Their results can be trivially extended to include the discount case.

The observable state case, that is the case where the true state of the system is known with certainty has been treated extensively. For both the finite and infinite horizon under a discount factor, Howard (1) developed a policy improvement routine for determining an optimal action and the optimal cost for each state.

II. Formulation

In this formulation, the notation of Smallwood and Sondik will be used. It is assumed that this system can be modeled by an N-state discrete time Markov decision process.

The observed state of the system is characterized by a probability vector π where π_i is the probability the true state of the system is i .

At each point in time an action must be selected from a finite set. Associated with an action, a , is a probability transition matrix P^a where P_{ij}^a is the conditional probability the system will make its next transition to state j given the current state is i and action a is taken. An observed outcome follows each action with $r_{j\theta}^a$ denoting the probability of observing output θ given the new state of the system is j and action a was taken. In addition an immediate reward $w_{ij\theta}^a$ is incurred if action a is taken, output θ is observed, and the system makes the transition from state i to state j . Thus if action a is taken and output θ is observed, the new state is π' where

$$\pi'_j = \left[\sum_i \pi_i P_{ij}^a r_{j\theta}^a \right] / \left[\sum_{ij} \pi_i P_{ij}^a r_{j\theta}^a \right] \quad (1)$$

The above transformation is summarized by

$$\pi' = T(\pi/a, \theta) \quad (2)$$

A policy is a rule that assigns an action to each possible state vector. It is required to find a policy that maximizes the expected discounted rewards over all periods for each possible state vector. Let $V(\pi)$ be the total discounted reward associated with such a policy.

Then $V(\pi)$ must satisfy the following recursive equation.

$$V(\pi) = \max_a \left[\sum_{i=1}^N \pi_i \sum_{j=1}^N p_{ij}^a \sum_{\theta} r_{j\theta}^a \left\{ w_{ij\theta}^a + \beta T(\pi/a, \theta) \right\} \right] \quad (3)$$

$$\text{Letting } q_i^a = \sum_{j,\theta} p_{ij}^a r_{j\theta}^a w_{ij\theta}^a \quad (4)$$

equation (3) is simplified somewhat to equation (5)

$$V(\pi) = \max_a \left[\sum_i \pi_i q_i^a + \beta \sum_{i,j,\theta} \pi_i p_{ij}^a r_{j\theta}^a V[T(\pi/a, \theta)] \right] \quad (5)$$

Once the function for $V(\pi)$ is known, an optimal action for π can be determined as one which maximizes the right hand side of (5).

III. A Learning Example

As an illustration, it will be shown how the system described in the previous section may be applied to the human learning process.

Consider a course which is given in several levels of instruction. The levels are denoted $1, 2, \dots, N$ with N being the easiest and 1 the hardest. The structure of the levels is a definite hierarchy in the sense that if a student knows the material at level i he must also know the material at any level $j > i$. Several examples where this situation may apply follow:

The first situation is one where the material covered at one level includes all that covered at preceding levels, plus some additional material. An example of this is a program developed at Behavioral Technology Laboratories (BTL) to teach students Kirchoff's Laws. This course is comprised of eleven levels with the lowest level defining the units for voltage, current and resistance up to the highest level which deals with the application of Ohm's Law and Kirchoff's voltage and current laws in complex networks. Another program developed at BTL is a short course in trigonometry consisting of five levels. At the lowest level students are given the definitions of the six basic trigonometric ratios. Then the student is given a right triangle in which the lengths of the sides are determined by a random number generator and the student is asked to determine these ratios for one of the acute angles. Succeeding levels deal with material on relationships between these ratios and problems testing the student's knowledge of these relationships.

A second situation is one where the material and problems covered at a particular level are virtually the same as the immediately preceding level except more clues and hints are given at the preceding level. A good example of this is a version of the Kirchoff's laws program considered

earlier at BTL in which problems would be given in level as follows:

1. Problems are given in steps with cues and knowledge of results at each step.
2. Problems are given in steps with no cues or knowledge of results at each step.
3. The student solves problems in steps but he chooses the steps.
4. The student is simply given problems and asked to solve them.

A third situation is one in which a student is to be drilled in a skill in order that he be able to perform it rapidly. Thus the exercises are virtually the same at all levels but the time constraints are tighter at the higher levels. In the BTL intercept trainer for the radar intercept observer function, the student is trying to fire a missile at the nose of a target and then turn around and fire another missile at the tail of that aircraft. The first missile is a radar guided missile fired when in the forward quarter and the second a heat seeker fired when in the rear quarter of the enemy aircraft. He is given a radar reading and must correct his angle of approach so as to be on a lead collision course that will insure a high hit probability when he fires the missile. At higher levels the student is given such problems at faster aircraft speeds.

Note, however, the assumption given for this model would not be applicable for the situation where a given level did not use certain material introduced at preceding levels.

A student is in state i if he knows the material of level i but not at any level more difficult than i and in state N+1 if he does not know the material at any level.

There are N actions and action i consists of instructing the student in the material of level i and then giving the student a test

on that material. For each action there are two possible outcomes-- either the student passes the test or he fails it. The objective is to develop an adaptive instructional sequence so that the student demonstrates knowledge of the material at level 1 as quickly as possible. Knowledge at level 1 is demonstrated by passing a test on the material at level 1. The reward, $w_{ij\theta}^a$, would be the negative of the expected time it would take to obtain instruction at level a and the system goes from state i to state j and θ (success or failure at a) is observed. For completeness a trap state ϕ would be needed. The student goes to state ϕ with probability one once he successfully completes the material at level 1. The only action in state ϕ is to do nothing which yields a zero reward and keeps the student in state ϕ with probability one.

Wollmer (6) treats the more restricted problem where $p_{ij}^a = 0$ unless $i=j$ or if $i=a$ and $j=i+1$. Thus if a student is in state i, he remains in state i unless he receives instruction at level $i+1$, in which case he either remains in state i or advances to state $i+1$. This would not allow the possibility of forgetting.

Other situations where partially observable Markov Decision processes occur are in machine replacement, decoding from sources transmitting over a noisy channel, medical diagnosis, and searching for a moving object.

Note, that if the assumption of a strict hierarchy in levels were dropped, the set of states would expand from $N+2$ to 2^N+1 including the trap state.

IV. The Maximum Reward Function

In this section it will be shown that a maximum reward function exists and that it is a convex function of the reward π .

Let $V_n(\pi)$ be the maximum reward function for the n period horizon. Then

$$V_n(\pi) = \max_a \left[\sum_i \pi_i q_i^a + \beta \sum_{i,j,\theta} p_{ij}^a r_j^\theta V_{n-1}[T(\pi/a, \theta)] \right] \quad (6)$$

Smallwood and Sondik (4) have shown that $V_n(\pi)$ is *

1. Convex

2. Piecewise Linear

$\lim_{n \rightarrow \infty} V_n(\pi)$ exists and is convex in π .

Define f_n so that $|V_n(\pi) - V_{n-1}(\pi)| \leq f_n$ all n and f_n is the smallest real number with this property and $V_0(\pi) = 0$. The f_n 's are well defined since all $V_n(\pi)$ are bounded above and below.

Lemma 1: $f_{n+1} \leq \beta f_n$

Proof : Choose $a(\pi)$ as the action that maximizes the right hand side of (6) for $V_{n+1}(\pi)$ if $V_{n+1}(\pi) \geq V_n(\pi)$ or for $V_n(\pi)$ otherwise.

Then $|V_{n+1}(\pi) - V_n(\pi)| \leq \left| \beta \sum_{i,j,\theta} p_{ij}^a r_j^\theta (V_n[T(\pi/a, 0)]) - V_{n-1}[T(\pi/a, 0)] \right| \leq \beta f_n$

$$|V_{n-1}[T(\pi/a, 0)]| \leq \beta f_n$$

Corollary 1: For $n^* > n$, $|V_{n^*}(\pi) - V_n(\pi)| < \epsilon(n)$

where $\epsilon(n) \rightarrow 0$.

* While Smallwood and Sondik assume $\beta=1$, their results hold for $0 < \beta \leq 1$.

Proof: From lemma 1, $f_n \leq \beta^{n-1} f_1$ and consequently

$$|v_n^*(\pi) - v_n(\pi)| \leq \sum_{i=n+1}^{n^*} f_i \leq \beta^n f_1 \sum_{i=0}^{\infty} \beta^i = f_1 \beta^n / (1-\beta)$$

Theorem 1: The function $v_n(\pi)$ is absolutely convergent.

Proof: Choose any particular $\pi = \bar{\pi}$. By Corollary 1, the $v_n^*(\pi)$ is bounded above and below and hence has an infinite convergent subsequence with limit $v^*(\pi)$. Choose $\epsilon > 0$ and n such that $\epsilon(n) < \epsilon$ for $N \geq n$ and $\epsilon(n)$ is as defined in corollary 1. For any $N \geq n$ and $\bar{n} \geq n$ in the convergent subsequence $|v_N(\pi) - v_{\bar{n}}(\pi)| < \epsilon$ and consequently $|v_N(\pi) - v^*(\pi)| < \epsilon$. Since n is independent of π , the theorem is proven.

Thus $V(\pi) = \lim_{n \rightarrow \infty} v_n(\pi)$ is well defined.

Theorem 2: $V(\pi)$ is convex in π .

Proof: Define $f(V, \pi_1, \pi_2) = V(\frac{1}{2}\pi_1 + \frac{1}{2}\pi_2) - \frac{1}{2}V(\pi_1) - \frac{1}{2}V(\pi_2)$.

Assume $V(\pi)$ is not convex and choose π_1 and π_2 such that $f(V, \pi_1, \pi_2) = k > 0$. Choose n such that $N > n \Rightarrow |v_N(\pi) - V(\pi)| < K/2$. $|f(v_N, \pi_1, \pi_2) - f(V, \pi_1, \pi_2)| < K$. Thus $f(v_N, \pi_1, \pi_2) > 0$ which is impossible since $v_N(\pi)$ is convex.

Note, that the piecewise linear property of $v_n(\pi)$ does not imply piecewise linearity of $V(\pi)$ as any continuous function may be expressed as the limit of a sequence of piecewise linear functions.

V. Linear Program Formulation

In the case of the observable finite state Markov decision processes with a discount factor, the problem of finding a maximum return for each state may be formulated as a linear program. The development of this may be found in Ross (6). In this section it is shown that a modification of this formulation extends to the problem formulated in Section II. Portions of the development which are similar to the finite state case will be outlined but without rigorous proofs.

Consider the set B of all continuous bounded functions defined on $S = \left\{ \pi / \pi_i \geq 0 \text{ all } i, \sum_{i=1}^n \pi_i = 1 \right\}$. Let the operator A be defined on this set as follows.

$$Au(\pi) = \max_a \left[\sum_i \pi_i q_i^a + \beta \sum_{i,j,\theta} \pi_i p_{ij}^a r_{j\theta}^a U[T(\pi)/a, \theta] \right] \quad (7)$$

Note that

1. $u \leq v \rightarrow Au \leq Av$
2. $Au \in B$ all $u \in B$
3. $A: B \rightarrow B$ is a contraction mapping on B .

The Operator A is the optimal return function for the one period problem in which a terminal reward $u(\pi)$ is given for the terminal state.

Since $A: B \rightarrow B$ is a contraction mapping, it has a unique fixed point $v = Av = \lim_{n \rightarrow \infty} A^n u$ for any $u \in B$. By Equation (3), this unique fixed point must be the optimal reward function. Let us consider any u such that $Au \leq u$. Then $u \geq Au \geq A^2 u \geq \lim_{n \rightarrow \infty} A^n u = v$. Thus the optimal return function v minimizes $u(\pi)$ for each $\pi \in S$ among all functions u satisfying $Au \leq u$.

In the finite state case where the above conditions also hold, it is noted that minimizing u_i for each state i may be accomplished by

minimizing the sum of the u_i 's. For this problem where such a sum would be infinite, the average value of $u(\pi)$ may be minimized. Thus, finding the function $u(\pi)$ is equivalent to solving the following infinite constrained program.

Find min Z , u such that

$$Z = \int \dots \int u(\pi) d\pi_n d\pi_{n-1} d\pi_{n-2} \dots d\pi_1 \quad (8)$$

$$\sum \pi_i = 1, \quad \pi_i \geq 0$$

subject to

$$\sum_i \pi_i q_i^a + \beta \sum_{i,j,\theta} \pi_i p_{ij}^a r_{j\theta}^a u[T(\pi/a, \theta)] \leq u(\pi) \text{ for } \quad (9)$$

$$\pi_i \geq 0, \quad \sum_i \pi_i = 1$$

Since the function $u(\pi)$ is continuous and defined on a closed bounded set, it may be expressed in an N -dimensional Maclaurin series:

$$V(\pi) = C_0 + \sum_{i_1, i_2, \dots, i_n} C_{i_1, i_2, \dots, i_n} \frac{\pi_1^{i_1} \pi_2^{i_2} \dots \pi_N^{i_N}}{i_1! i_2! \dots i_n!} \quad (10)$$

If $V(\pi)$ is expressed as such a series or approximated by a partial series consisting of terms up to degree n , the coefficient of C_{i_1, i_2, \dots, i_n} in (8) is simply

$$\int_{\pi_1=0}^1 \int_{\pi_2=0}^{1-\pi_1} \dots \int_{\pi_N=0}^{1-\pi_1-\pi_2-\dots-\pi_{N-1}} \frac{1}{i_1! i_2! \dots i_n!} d\pi_N d\pi_{N-1} \dots d\pi_1 \quad (11)$$

In evaluating the integral the following lemma is needed.

$$\text{Lemma 2: } \int_0^a (a-x)^m x^n dx = \frac{m! n!}{(m+n+1)!} a^{m+n+1}$$

Proof: Integrating by parts one obtains for the above

$$\text{integral} = \frac{1}{m+1} x^{n-1} (a-x)^m \Big|_0^a + \frac{n}{m+1} \int_0^a (a-x)^{m+1} x^{n-1} dx \\ = \frac{n}{m+1} \int_0^a (a-x)^{m+1} x^{n-1} dx. \text{ Applying this relationship}$$

$$\text{recursively, one obtains } \frac{n! m!}{(m+n)!} \int_0^a (a-x)^{m+n} dx = \frac{m! n!}{(m+n+1)!} a^{m+n+1}$$

From this lemma, expression (11) can be evaluated.

Theorem 3: The value of expression (11) is $\prod_{j=1}^n i_j! / \left[\sum_{j=1}^n (i_j + 1) \right]!$

Proof: Integrating (11) with respect to π_n gives

$$\frac{i_n!}{(i_n+1)!} \int_0^1 \pi_1^{i_1} \int_0^{1-\pi_1} \pi_2^{i_2} \cdots \int_0^{1-\sum_{j=1}^{n-1} \pi_j} (1 - \sum_{j=1}^{n-1} \pi_j)^{i_n+1} \pi_{n-1}^{i_{n-1}} d\pi_{n-1} \cdots d\pi_1$$

Applying lemma 1 with $a = 1 - \sum_{j=1}^{n-1} \pi_j$ and integrating with respect to π_{n-1} yields

$$\frac{i_n! i_{n-1}!}{(i_n+1 i_{n-1}+2)!} \int_0^1 \pi_1^{i_1} \int_0^{1-\pi_1} \pi_2^{i_2} \cdots \int_0^{1-\sum_{j=1}^{n-2} \pi_j} (1 - \sum_{j=1}^{n-2} \pi_j)^{i_n+1 i_{n-1}+2} d\pi_{n-2} \cdots d\pi_1$$

Continual application of lemma 2 yields $\prod_{j=1}^n i_j! / \left(\sum_{j=1}^n (i_j + 1) \right)!$

Thus if $V(\pi)$ is to be approximated by an n^{th} degree polynomial function in π , then substituting the expression of theorem 3 and (1) in (8) and (9) and rearranging terms yields:

Find $c_0, c_{i_1 i_2 \dots i_n}$, $\min z$ such that

$$z = c_0 + \sum \left[\left(\prod_{j=1}^n i_j! \right) / \left(\sum_{j=1}^n (i_j + 1)! \right) \right] c_{i_1 i_2 \dots i_n} \quad (12)$$

$$(1-\beta)c_0 + \sum \left(\prod_{j=1}^n \pi_j^{i_j} - \beta \sum_{\theta} [k_{i_1 i_2 \dots i_N}^{(\theta)} / d_{i_1 i_2 \dots i_N}^{(\theta)}] \right) \times$$

$$c_{i_1 i_2 \dots i_N} \geq \sum_{i=1}^n \pi_i q_i^a \quad (13)$$

$$\text{where } k_{i_1 i_2 \dots i_N}^{(\theta)} = \prod_{j=1}^n (\sum_{i_j} \pi_i p_{ij} r_{j\theta}^a)^{i_j} \quad (14)$$

$$d_{i_1 i_2 \dots i_N}^{(\theta)} = (\sum_{i_j} \pi_i p_{ij} r_{j\theta}^a)^{\sum (i_j - 1)} \quad (15)$$

for all θ , all $\pi \geq 0$ such that $\sum \pi_i = 1$

Thus the problem of solving the program (8-9) with a multinomial approximation of $u(\pi)$ becomes a linear program (12-15) with an infinite number of constraints and unrestricted variables. Note that the minimum value of Z obtained in the linear program (12-15) would actually be larger than that obtained in the program (8-9).

VI. Computational Procedure

Given an optimal solution to the linear program (12-15), consider the set of constraints for which the $C_{i_1 i_2 \dots i_N}$ are basic. If the program was solved with these constraints only, the same solution would be obtained and all other constraints would be satisfied. Thus, while the program consists of an infinite number of constraints, only a finite number need to be included provided the correct ones are chosen. This will be taken advantage of by solving the program with a finite subset of the constraints, introducing an unsatisfied constraint, then dropping any that are not binding, and continuing until an optimal solution is obtained.

Let the quantity $f(\pi, C)$ be defined as follows:

$$F(\pi, C) = (1-\beta)C_0 + \sum_{j=1}^n \pi_j^{i_j} - \beta \sum_{\theta} [k_{i_1 i_2 \dots i_N}^{(0)} d_{i_1 i_2 \dots i_N}^{(0)}]$$
$$C_{i_1 i_2 \dots i_N} = \sum_{i=1}^n \pi_i^{q_i^a} \quad (16)$$

The constraints (13) are equivalent to $F(\pi, C) \geq 0$ all π . Thus if at least one constraint is not satisfied for a given C vector, the value of π that minimizes $F(\pi, C)$ is the most unsatisfied one.

The procedure for solving the linear program (12-15) is given in algorithm 1.

Algorithm 1

1. Formulate the linear program with any finite subset of the constraints in (13).
2. Solve the linear program for C .
3. Delete any constraints for which a slack variable is basic.
4. Solve the following non-linear program.

Find $\pi \geq 0$, min Z such that

$$Z' = f(\pi, C) \quad (17)$$

$$\sum_{i=1}^N \pi_i = 1 \quad (18)$$

If $Z' \geq 0$, terminate as C is optimal. Otherwise introduce the constraint corresponding to the value of π that optimizes (17-18) and go back to Step 2.

A local optimum to (17-18) may be found by algorithm 2.

Algorithm 2

1. Choose an arbitrary probability vector and evaluate $f(\pi, C)$.
2. Find an order pair (i, j) such that increasing π_i by ϵ and decreasing π_j by ϵ decreases $f(\pi, C)$ without violating $0 \leq \pi_i \leq 1$ and $0 \leq \pi_j \leq 1$. If no such pair can be found, terminate as π is a local optimum.
3. Increase π_i to $\bar{\pi}_i$ and decrease π_j to $\bar{\pi}_j$ such that neither the pair (i, j) or (j, i) satisfied the conditions of Step 2. Then go back to Step 2.

For finiteness, the ϵ of Step 2 would be chosen ahead of time.

There are several ways of performing Step 3 to find the new value of π_i and π_j . One efficient way is to first bracket π_i and π_j between π'_i , π'_j and π''_i and π''_j and continually reduce the difference between these by a factor of one half, thus converging on a single point.

Initially π'_i and π'_j would be the current values of π_i and π_j and $\pi''_i = \pi_i + \delta$, $\pi''_j = \pi_j - \delta$ where $\delta = \min [1 - \pi_i, \pi_j]$. Then consider the pair $\bar{\pi}_i = \frac{1}{2}(\pi'_i + \pi''_i)$ and $\bar{\pi}_j = \frac{1}{2}(\pi'_j + \pi''_j)$. If $f(\bar{\pi}, C)$ is a local

minimum under the restriction that all components of π other than π_i and π_j are held constant, then $\bar{\pi}$ is the desired point. Otherwise, let $\bar{\pi}_i$ and $\bar{\pi}_j$ replace π_i and π_j if the direction of decrease is towards π_i and π_j but let $\bar{\pi}_i$ and $\bar{\pi}_j$ replace π_i and π_j if the direction of decrease is towards π_i and π_j . If neither direction yields a decrease, let $\bar{\pi}_i$ and $\bar{\pi}_j$ replace π_i and π_j if $f(\bar{\pi}) > f(\pi)$ but replace $\bar{\pi}_i$ and $\bar{\pi}_j$ otherwise. Step 3 would terminate when $\pi_i - \bar{\pi}_i < \epsilon_1$ where $\epsilon_1 < \epsilon$.

Note that if the C vector approximation of $U(\pi)$ were exact, any local minimum of $f(\pi, C)$ would be a global minimum due to the convexity of $V(\pi)$. While this is not guaranteed in the approximation, one could take random samples of π in an attempt to find a vector yielding a lower value of π than the local minimum or evaluate Z' for all π vectors whose components are multiples of $1/n$ where n is large if the result $\min Z' = 0$ is obtained.

When introducing an unsatisfied constraint, it is recommended that the dual simplex method be used to solve the resulting program which is already dual feasible.

The sequence of $\min Z$ values generated by algorithm 1 is non-decreasing, bounded above, and hence must have a limit. It is an open question as to whether this limit is the true $\min Z$ or in particular if the sequence of Z' values in algorithm 2 tend to zero. Consider the sequence of linear programs solved by algorithm 1 and assume the number of equations in each equals the number of components in the C vector plus one. It has already been shown that it will not exceed this number and if it is less, additional constraints with all coefficients being zero may be added. Consider also the sequence of matrices formed by the probability vectors that generate these constraints. Since these

are bounded above, these matrices, and consequently the set of linear programs for algorithm 1 must have a convergent subsequence. Consider now the sequence of constraints generated by this sequence in algorithm 2. By the same argument this sequence must have a convergent subsequence. In this latter sequence, either $f(\pi, C) \rightarrow 0$ or else the cost coefficient in the pivot column tends to zero for if not the increase in $\min z$ would not tend to zero which is impossible since $\min z$ is bounded above.

If the sequence of $f(\pi, C)$ values generated by problem 2 did not appear to tend to zero after many iterations while the change in $\min z$ did appear to tend to zero, some possible ways out are as follows. First one may sample a large number of probability vectors and find one which would give the largest increase in z on a single pivot. Second, one may search all probability vectors that are multiples of $1/n$ where n is a large number and find the one which gives the largest increase in z for one pivot.

It should be noted that if the sequence of z values obtained in algorithm 2 do not tend to zero, then one has a situation somewhat analogous to cycling in the dual simplex method. Since cycling almost never occurs in the primal simplex method, there appears to be some basis for thinking that the sequence of z values would tend to zero the majority of times.

One could of course only consider constraints generated by probability vectors whose components are multiples of $1/n$. By imposing a lexicographic ordering, one could insure a true optimum in a finite number of steps.

VII. Bounds on Accuracy

In solving the non-linear program (17-18) in Step 4 of the algorithm to find the most unsatisfied constraint of the linear program (12-15), one may wish to terminate the program when $Z > \delta$ rather than for $Z \geq 0$ where δ is a small positive number. If so, the value of Z obtained for (12) will be less than the true minimum for Z , since the program has been optimized for only a subset of the constraints. However, it is easy to see from (12) and (13) that increasing C_0 by $\delta/(1-\beta)$ yields a feasible solution and increases Z by that same amount. Consequently, this feasible set would come to within $\delta/(1-\beta)$ of minimizing Z .

The question now arises as to how close $V(\pi)$, the Maclaurin series approximation to $V(\pi)$, is to the true value of $V(\pi)$. To answer this consider the operator $Au(\pi)$ defined in equation (7) and define:

$$\max_{||Au - u||} = \pi ||Au - u|| \quad (19)$$

Since the operator A is a contraction mapping with $|Au - Av| \leq \beta|u - v|$ it can be shown that $||A^{n+1}u - A^n u|| \leq \beta^n ||Au - u||$ and $||A^n u - u|| \leq (1-\beta^n) ||A^n u - u|| / (1-\beta)$ and $V(\pi) = \lim_{n \rightarrow \infty} A^n u$, it follows that

$$|V(\pi) - \tilde{V}(\pi)| \leq |Av - v| / (1-\beta) \quad (20)$$

One could find a local maximum to $|Av - v|$ by an incremental procedure similar to that used to find the most unsatisfied constraint to introduce into the linear programming problem. Alternatively, one could enumerate (20) for all possible probability vectors whose components are multiples of $1/n$.

REFERENCES

1. Howard, R. A. Dynamic Programming and Markov Processes, John Wiley and Sons, New York, 1960.
2. Manne, A. "Linear Programming and Sequential Decisions," Management Science 6, 259-267 (1960).
3. Ross, S. M. Applied Probability Models with Optimization Applications, Holden-Day, San Francisco, 1970.
4. Smallwood, R. D. & E. J. Sondik, "The Optimal Control of Partially Observable Markov Processes Over a Finite Horizon," Operations Research 21, 1071-1087, (1973).
5. Wolfe, P., and Dantzig, G. B. "Linear Programming in a Markov Chain, Operations Research 10, 702-710 (1962).
6. Wollmer, R. D., "A Markov Decision Model for Computer-Aided Instruction," Behavioral Technology Laboratories, University of Southern California, Technical Report #72, December 1973.