A THEORY FOR SEMI-MARKOV DECISION PROCESSES WITH UNBOUNDED COSTS AND ITS APPLICATION TO THE OPTIMAL CONTROL OF QUEUEING SYSTEMS 1

 \mathbb{R}_{p}

5.

AD A () 3 () 6 4



TECHNICAL REPORTENO. 64 AUGUST 1976

PREPARED UNDER CONTRACT NC0014-76-C-0418 (NR-047-061) FOR THE OFFICE OF NAVAL RESEARCH

Reproduction in Whole or in Part is Permitted for any Purpose of the United States Government This document has been approved for public release and sale; its distribution is unlimited

> DEPARTMENT OF OPERATIONS RESEARCH STANFORD UNIVERSITY STANFORD, CALIFORNIA



11

PETER ORKENYI

1976

DET 13

A THEORY FOR SEMI-MARKOV DECISION PROCESSES WITH UNBOUNDED COSTS, AND ITS APPLICATION TO THE OPTIMAL CONTROL OF QUEUEING SYSTEMS

by

PETER ORKENYI

TECHNICAL REPORT NO. 64 AUGUST 1976

PREPARED UNDER CONTRACT N00014-76-C-0418 V (NR-047-061) FOR THE OFFICE OF NAVAL RESEARCH

Frederick S. Hillier, Project Director

Reproduction in Whole or in Part is Permitted for any Purpose of the United States Government

This document has been approved for public release and sale; its distribution is unlimited.

This research was supported in part by NATIONAL SCIENCE FOUNDATION GRANT ENG 75-14847

DEPARTMENT OF OPERATIONS RESEARCH STANFORD UNIVERSITY STANFORD, CALIFORNIA

100:3 104 61 171s With C. Ching Ε γ U.A. 2002) I del Sector o · 103110,1 \square land artical Addition

CHAPTER 1

INTRODUCTION

Markov and semi-Markov decision processes have been studied extensively since their initial developmant in the late 1950's and early 1960's. They provide the natural framework for the study of a plethora of problems arising in the areas of queueing, inventory, maintenance and replacement, etc. Many useful results about Markov and semi-Markov decision processes are available now under a variety of assumptions. A common assumption has been the assumption of bounded costs. Although bounded costs is an appropriate assumption for many problems, there are also many situations, especially in the context of queueing and inventory, for which it is not appropriate. Thus, there is a need for developing a unity for Markov and semi-Markov decision processes with unbounded costs. Although there have been some efforts in this direction ear 'ier, stronger results need to be developed. That is the objective of this report. Specifically, results are obtained for semi-Markov decision processes both when the costs are discounted _nd when they are not. Application to the optimal control of queueing systems is also considered.

The terminology of semi-Markov decision processes is summarized in Section 1. Section 2 then presents some examples of semi-Markov decision processes both with and without unbounded costs. Section 3 reviews the literature on semi-Markov decision processes. An overview of the study is presented in Section 4.

1

a said a subolight to the address of the said of

1. Terminology of Semi-Markov Decision Processes.

The semi-Markov decision process is a stochastic process which requires certain decisions to be made at certain points in time. These points in time are the <u>decision epochs</u>. At each decision epoch, the system under consideration is observed and found to be in a certain <u>state</u>. The set of all conceivable states is <u>the state space</u>. The decision consists of choosing an <u>action</u> from a set of permissible actions. This set depends on the state of the system when the decision has to be made. The set of permissible actions for a given state is an <u>action space</u>. The union of all action spaces is referred to as <u>the action space</u>. Once an action has been chosen, the probabilistic aspects of the evolution of the system until the next decision epoch occurs (including the time elapsed and the state of the system at the next decision epoch) is completely determined by the state of the system when the action was chosen and the action itself.

A <u>policy</u> for a semi-Markov decision process is a rule which selects an action at each decision epoch by considering only the history of the process up to that point in time. An interesting class of policies is the class of <u>stationary</u> policies. A stationary policy selects the action at each decision epoch solely on the basis of the state of the system at the decision epoch. A stationary policy is <u>deterministic</u> if it selects the actions according to a fixed mapping from the state space into the action space; otherwise it is randomized.

A part of the process is the costs incurred. The objective is to minimize these costs. They are, however, incurred in a random fashion and at different times, so a further specification of the objective is needed. There are several alternatives. If the time factor is not

important, one may choose to minimize the total expected cost, or if this is not finite, the long-run expected average cost. If the time factor is important, one may discount the costs and minimize the total expected discounted cost.

For our purposes, a semi-Markov decision process is completely specified by four objects, the state space, S, the action spaces $\{A_{\alpha}\}_{\alpha \in \Omega}$, the law of motion q, and the cost function c. Let $A = \bigcup_{s \in S} A_s$, and let R be the set of real numbers. The law of motion, q, is a mapping from $S \times A \times S \times R$ into R, and the cost function, c, is a mapping from $S \times A \times R$ into R. Consider a decision epoch. Suppose the state there is s and suppose the action chosen there is a. Then, for s' \in S and t \in R, q(s,a,s',t) is the joint probability that the time until the next decision epoch is less than or equal to t and that the state at the next decision epoch is s'. If the times between the decision epochs are constant, then we have a Markov decision process. Also, for $t \in R$, c(s,a,t) is the expected cost accumulated until time t. The formulation of a problem in the framework of semi-Markov decision processes consists of specifying S, ${A_s}_{s \in S}$, q and c. Some examples of semi-Markov decision processes are now presented.

2. <u>Examples of Semi-Markov Decision Processes With and</u> .thout Unbounded <u>Costs</u>.

Selling an asset (Ross (1970)):

Consider a person who wants to sell his house. Offers arrive according to a stationary Poisson process. The sizes of the offers are independent, identically distributed random variables. When an offer arrives, it

must either be accepted or rejected. Rejected offers are lost. A maintenance cost is incurred at a constant non-negative rate until the house is sold. The problem is to decide when an offer should be accepted. This problem can be formulated within the framework of a semi-Markov decision process as f lows.

Let the accision epochs be the same as the epochs when offers arrive, let the actions be to accept or reject the current offer, and let the state of the system be the size of the offer at the most recent decision epoch.

A job shop model (Lippman and Ross (1968)):

Consider a factory which is only able to handle one job at a time. Jobs arrive according to a stationary Poisson process. When a job arrives it is classified to be of a certain type. Jobs of the same type have an identical probabilistic structure for their cost and completion time. The classification of arriving jobs are independent, identically distributed random variables. Each job must either be accepted or rejected. Jobs arriving when the factory is busy are rejected automatically. The problem is to determine when a job should be accepted (rejected) when the factory is not <u>busy</u>. This problem can be formulated within the framework of semi-Markov decision processes as follows.

Let the decision epochs be the same as the epochs of job arrivals (neglect jobs which arrive when the factory is busy), let the available actions be to accept or reject the job that just arrived, and let the state of the system be the type of job present.

The M/G/l queueing system with removable server (Heyman (1968)): Consider a queueing system having one server which can be turned

on and off. Customers arrive according to a stationary Poisson process. They are served one by one on a first-come-first-served basis. The service times are independent, identically distributed random variables. There is a cost associated with the service of each customer. These costs are independent, identically distributed random variables. There are fixed charges for turning the server on and off. There is a cost for having the server on when there are no customers in the system. This cost is incurred at a constant rate at such times. Finally, there is a cost for holding customers in the system. This cost is incurred at a rate which is a non-negative, non-decreasing function of the number of customers present. The problem is to determine when the server should be turned on and turned off. This problem can be formulated within the framework of semi-Markov decision processes as follows.

Let the decidion epochs be the epochs of customer arrivals and departures (neglect arrivals which occur when the server is busy). Let the available actions be to turn the server off (or have him off) and to turn him on (or have him on). Finally, let the state of the system be a vector whose first component gives the number of customers present, and whose second component shows the status of the server.

3. A Brief Survey of the Literature on Semi-Markov Decision Processes.

The first comprehensive study of Markov decision processes was done by Howard (1960). He assumed finite state and action spaces, and considered the problem both with and without discounting. He only considered stationary policies, and developed his now well-known policy improvement procedures. He proved that they would produce optimal stationary policies.

At the same time, Manne (1960) suggested solving the Markov decision

problem by using linear programming. He used the average cost criterion, and showed how to solve an inventory problem by his suggested approach. The first linear programming formulation for he problem with discounting was given by d'Epenoux (1960). Shortly afterwards, Wolfe and Dantzig (1962) proposed the use of their decomposition technique on Manne's linear programming formulation.

Blackwell (1962) considered Markov decision processes with finite state and action spaces, and proved that there is a stationary policy which is optimal among all Markov policies. He also considered the problem for arbitrarily small interest rates, and proved that there is a stationary policy which is optimal among all Markov policies for small enough interest rates. Later, Blackwell (1965) considered Markov decision processes with more general state and action spaces. He only assumed that they were Borel sets. However, he assumed that the rewards were uniformly bounded. He considered the problem with discounting, and allowed any measurable policy. His main results were the following. There is a (p,ϵ) -optimal stationary policy. If the action spaces are countable, then there is an ϵ -optimal stationary policy. If the action spaces are finite, then there is an optimal stationary policy. If there is an optimal policy, then there is one which is stationary.

Strauch (1966) considered the same problem as Blackwell, but instead of using discounting, he assumed that the rewards were negative. His main results were similar to those of Blackwell. If the action spaces are finite, then there is an optimal policy. If there is an optimal policy, then there is one which is stationary. The optimal return function is measurable and satisfies the optimality equation.

Denardo (1967) also considered the same problem as Blackwell and generalized it to include certain stochastic games. He introduced operators with certain monotonicity and contraction properties, and used the Picard-Banach fixed point theorem to prove that the functional equation of optimality has a unique solution, which is the optimal reward function.

Veinott (1966) gave a policy iteration procedure for finding a biasoptimal policy (no discounting). Later, Veinott (1969) considered a more refined optimality criterion, namely, that of finding a policy which is optimal for all sufficiently small interest rates (sensitive discount optimality). He developed a policy iteration procedure for finding a stationary policy which would be optimal according to this criterion.

Derman (1966) considered Markov decision processes with finite action spaces and a countable state space. He used the average cost criterion, and gave conditions for when a stationary, deterministic policy is optimal. Ross (1968) considered the same problem, but allowed a general state space. He derived results similar to those of Derman. He also suggested a method for converting the average cost problem to a discounted cost problem.

One of the first to consider semi-Markov processes was Pyke (1961). Shortly afterwards, Howard's results for Markov decision processes were extended to semi-Markov decision processes independently by Jewell (1963) and Howard (1964). When they considered the average cost criterion, they assumed that all states belong to one positive recurrent class. They also gave linear programming formulations.

Denardo and Fox (1968) considered the multi-chain case (i.e., the case of several positive recurrent classes), using the average cost criterion. They gave a linear programming formulation and a policy

improvement procedure. Later Denardo (1970a) developed a solution method which used Manne's linear programming formulation to solve a sequence of subproblems. This solution method has the advantage that several small linear programming problems are solved instead of one big one. Denardo (1971) also considered the problem when small interest rates are used. His results are similar to those of Veinott for the discrete-time Markov decision process. He gives a sequence of linear programming problems for finding an optimal policy.

All of these authors have assumed that the immediate rewards or costs are bounded uniformly. After Strauch, Harrison (1972) was the first one to relax the condition of bounded costs. He assumed that the expected absolute reward in one period minus the expected absolute reward in the period before it, given the state at the beginning of that period, is uniformly bounded. He then showed that the expected discounted reward is finite for each policy and that there exists a stationary policy which is optimal. He proved this by using the Picard-Banach fixed point theorem. He also extended his results from Markov decision processes to semi-Markov decision processes.

The problem with unbounded costs was also considered by Reed (1973). He investigated the problem both with and without discounting. He assumed finite action spaces and countable state space. He gave sufficient conditions for a stationary policy to be optimal.

Hordijk (1974a), (1974b) also considered the problem with unbounded costs. He introduced the notion of <u>convergent</u> dynamic programming, which is just to say that the expectation of the sum of the absolute rewards is finite. He proved that a policy is optimal if it is unimprovable and if another condition is satisfied.

Most recently, Lippman (1973), (1975a) considered the problem with unbounded costs. His approach is to use a norm such that the norm of the costs is finite even though the costs are unbounded. In order to obtain the usual results, he then has to make ascumptions about the law of motion of the system. By doing that, he showed that Denardo's N-stage contraction assumption is satisfied, and the results follow.

4. Overview of the Study.

The emphasis of this report is on *i*_cervining necessary and sufficient conditions for a stationary policy to be optimal. It is not assumed that the costs are bounded. The problem is considered both with and without discounting.

Chapter 2 treats the problem without discounting. Two closely related optimality criteria are used, namely, the average cost criterion and the undiscounted cost criterion. After introducing the important concept of an unimprovable policy, sufficient conditions are given f in unimprovable policy to be optimal. Both the special case where the optimal expected average cost is independent of the start-state and the general case when the average cost is not necessarily constant are considered.

Chapter 3 treats the r is m with discounting. After formulating the problem and introducing we operators Q_T and T_T , the optimality equation is proven. The existence of stationary optimal and stationary e-optimal policies are then investigated. Policy improvement is considered, and some necessary and sufficient conditions for optimality are given.

Chapter 4 is devoted to the optimal control of queueing systems.

Solution methods are explored, and four different ways of solving the problem of unbounded costs are presented.

Some general notation and conventions are best introduced here. R denotes the set of real numbers, R_+ denotes the set of non-negative real numbers, N denotes the set of natural numbers (starting with one) and N_0 denotes the non-negative integers. The Kroene-ner delta function δ is defined by

$$\delta(x,y) = \begin{cases} 1 & \text{if } x = y , \\ 0 & \text{if } x \neq y . \end{cases}$$

If x is a real number, then x^+ is max(0,x) and x^- is max(0,-x). Finally, we use the convention that

$$x + y = \begin{cases} \infty & \text{if } x = \infty, y > -\infty, \\ -\infty & \text{if } x < \infty, y = -\infty, \\ \text{undefined if } x = -y = +\infty. \end{cases}$$

CHAPTER 2

SEMI-MARKOV DECISION PROCESSES WITHOUT DISCOUNTING

This chapter presents an investigation of semi-Markov decision processes without discounting the costs. Thus, costs of equal size incurred at different times count the same. Two optimality criteria are used. The first one is the average cost criterion, according to which a policy is optimal if the long-run expected average cost is minimized by this policy. This criterion has been considered recently by Hordijk (1974a). The other criterion is the <u>undiscounted cost</u> <u>criterion</u>. A policy is optimal under this criterion if it minimizes the long-run (total) expected cost for the process which is derived from the original one by incurring an additional cost at a rate equal to the negative of the minimum average cost. This criterion has been considered by Denardo (1970). He called a policy which is optimal for this criterion a bias-optimal policy.

There have traditionally been two approaches to the problem without discounting. The first one consists of restricting one's consideration to stationary (deterministic) policies and performing a stationary analysis. The second one consists of considering the problem with discounting and observing what happens when the interest rate goes to zero. Here, we will follow the first approach. It has been common to assume that the costs are uniformly bounded. We make no assumptions about the size of the costs. Reed (1973) conducted a similar but somewhat less complete study of the problem.

In Section 1, there is a formal statement of the problem to be considered. It also contains some preliminary results. <u>Unimprovable policies</u> are defined there. In Section 2, sufficient conditions for an unimprovable policy to be optimal are given. It is assumed that the long run expected average cost is constant. In Section 3, the results from Section 2 are extended to cover the general case of non-constant long-run expected average cost. In Section 4, there is a brief discussion of methods for finding an optimal policy.

1. Problem Formulation.

As before, let S be the state space, $(A_s)_{s\in S}$ be the action spaces, q be the law of motion and c the cost function. Let O be the set of stationary, deterministic policies, and let A be $\bigcup_{s\in S} A_s$. For each $n \in N$, let t_n , s_n , and a_n denote the time of the n^{th} decision epoch, the state observed there, and the action chosen there, respectively.

For each $\pi \in \mathfrak{S}$, let v_{π} be the mapping from $\mathbb{S} \times \mathbb{R}^+$ into \mathbb{R} such that, for each $s \in \mathbb{S}$ and $t \in \mathbb{R}^+$,

$$v_{\pi}(s,t) = E_{\pi,s}\left\{\sum_{n \in N_{t}} c(s_{n},a_{n},t-t_{n})\right\},$$

where

$$N_t = \{n \in N | t_n \leq t\}$$
.

E is the expectation operator, and the subscripts π and s respectively denote that the start-state is s and that the policy used is π . In words, $v_{\pi}(s,t)$ is the expected cost incurred until time t, given that the start-state is s and that the policy π is used.

 v_{π} need not always be well-defined. Later, however, certain assumptions which guarantee the existence of v_{π} for each $\pi \in \mathcal{D}$ will be made.

The analysis here is based on the fact that under certain conditions (to be introduced when needed), $v_{\pi}(t,s)$ has a linear asymtote for each $s \in S$ and $\pi \in D$. For each $\pi \in D$, let φ_{π} and w_{π} be the mappings from S into R such that

$$\begin{split} \varphi_{\pi}(s) &= \lim_{t \to \infty} v_{\pi}(s,t)/t , \\ w_{\pi}(s) &= \lim_{t \to \infty} \{v_{\pi}(s,t) - t \cdot \varphi_{\pi}(s)\} \end{split}$$

for $s \in S$. φ_{π} is the long-run expected average cost, given that the start state is s and that the policy π is used. $w_{\pi}(s)$ is the long-run expected cost not accounted for by $\varphi_{\pi}(s)$.

Two optimality criteria will be used. The first one is the <u>average</u> <u>cost criterion</u>. A policy $\pi^* \in \mathcal{O}$ is optimal according to this criterion if $\varphi_{\pi}(s) \leq \varphi_{\pi}(s)$ for $s \in S$ and $\pi \in \mathcal{O}$, and the policy is called <u>average optimal</u>. The second criterion is the <u>undiscounted cost criterion</u>. A policy $\pi^* \in \mathcal{O}$ is optimal according to this criterion if it is average optimal and, in addition, $w_{\pi}(s) \leq w_{\pi}(s)$ for $s \in S$ such that $\varphi_{\pi}(s) = \varphi_{\pi}(s)$ for $\pi \in \mathcal{O}$. A policy which is optimal in this sense is called <u>undiscounted optimal</u>. This latter criterion has not received much attention in the literature. This may be due to the fact that often there is not much to gain by using this criterion instead of the average cost criteria is that the action in the transient states become more important when the undiscounted cost criterion is used. To illustrate this point further, an example is included below.

Example: Consider the following simple semi-Markov decision process. The state space is N and the action spaces are {0,1}. The times between the decision epochs are exponentially distributed with the same parameter. State 0 is an absorbing state. Consider states in N. If action 0 is taken, the state 0 is entered next with probability one. If action 1 is taken, the state numbered 1 higher is entered next with probability one. The cost structure is simple. Each time a state in N is reached, an immediate cost of 2 units is incurred, and each time the state 0 is entered, an immediate cost of 1 unit is incurred. Any policy which chooses action 0 in all the states above a given number is average optimal. The undiscounted optimal policy is the one which always chooses action 1. This is clearly the desired policy.

One special reason for using the undiscounted cost criterion is as follows. Under certain circumstances there may exist a sequence of average optimal policies π_1, π_2, \ldots such that using π_1 for the first decision, π_2 for the second, π_3 for the third, and so on, leads to a long-run expected average cost which is higher than the optimal one. This can easily be seen from the example above. First let π_n be the policy which chooses action 1 for states numbered less than n and action 0 for states numbered n or higher. Each π_n is average optimal. But using π_n at the nth decision epoch for $n = 1, 2, \ldots$, leads to a long-run expected average cost twice as high as the optimal one. Notice that since there is a unique undiscounted optimal policy, this situation cannot occur when the undiscounted cost criterion is used. In general, there is no guarantee for the existence of a unique undiscounted optimal policy, but often a unique undiscounted optimal policy does exist

and thus the undesirable situation mentioned above can be avoided by using the more refined criterion. Some useful semi-Markov process terminology will now be introduced.

A state is called <u>transient</u> if with probability one it will not be reentered after some time. A state is called <u>recurrent</u> if with probability one it will always be reentered. A recurrent state is <u>positive recurrent</u> if the expected time between consecutive visits of this state is finite. Otherwise, it is called <u>negative recurrent</u>. If there is a positive probability that a state is reached in a finite time from another state and vice versa, then the two states are said to <u>communicate</u>. The positive recurrent states belong to one or more <u>positive</u> <u>recurrent classes</u> of states. Each positive recurrent class is a set of positive recurrent states which communicate with each other, but not with states outside the class. We make the following assumptions.

Assumption 1: There is an $\epsilon > 0$ such that

「たいで、たいである

ないないであったとうないできょうできょうというないであったとうない

 $q(s,a,s',\epsilon) = 0$, for $s \in S$, $a \in A_s$, $s' \in S$.

In words, the time between two consecutive decision epochs is at least ϵ .

<u>Assumption 2</u>: For each $\pi \in \mathcal{D}$ and $s \in S$, the expected cost incurred and the expected time elapsed from time t until the first decision epoch after (or at) time t divided by the time t have zero as their limits as t tends to infinity, given the start-state s and policy π .

Faced with a particular semi-Markov decision process, one may have difficulties in showing that it satisfies the above assumption. However,

we have not been able to do without them. If the semi-Markov decision process is a Markov decision process, then the second assumption is trivially satisfied.

Some convenient notation will now be introduced. For each $\pi \in \mathcal{D}$, let q_{π} and τ_{π} be the mappings from S × S into R such that

$$q_{\pi}(s,s') = \lim_{t \to \infty} q(s,a_{\pi}(s),s',t) ,$$

$$\tau_{\pi}(s,s') = \int_{t \in R_{+}} tdq(s,a_{\pi}(s),s',t) ,$$

for s,s' \in S. $a_{\pi}(s)$ is the action chosen by π in the state s. For each $\pi \in \mathcal{D}$, also let v_{π} and c_{π} be the mappings from S into R such that

$$v_{\pi}(s) = \sum_{s' \in S} \tau_{\pi}(s, s') ,$$

$$c_{\pi}(s) = \lim_{t \to \infty} c(s, a_{\pi}(s), t) ,$$

for $s \in S$. $q_{\pi}(s,s^{\circ})$ is the probability that the next state will be s, given the present state s and policy π . $\tau_{\pi}(s,s^{\circ})$ is $q_{\pi}(s,s^{\circ})$ multiplied by the expected time until the next decision epoch, given that the next state is s'. $\tau_{\pi}(s)$ is the expected time until the next decision epoch, given the present state s and policy π . $c_{\pi}(s)$ is the expected cost until the next decision epoch, given the present state s and policy π . Naturally, we assume that all these quantities exist and are finite.

$$\begin{split} & \varphi_{\pi}(s) = \sum_{s^{\dagger} \in S} q_{\pi}(s,s^{\dagger}) \cdot \varphi_{\pi}(s^{\dagger}) , \\ & w_{\pi}(s) = c_{\pi}(s) - \sum_{s^{\dagger} \in S} \tau_{\pi}(s,s^{\dagger}) \cdot \varphi_{\pi}(s^{\dagger}) + \sum_{s^{\dagger} \in S} q_{\pi}(s,s^{\dagger}) \cdot w_{\pi}(s^{\dagger}) , \end{split}$$

for $s \in S$ and $\pi \in \mathcal{D}$ (see Denardo and Fox (1968)). The expressions on the right-hand side are obtained by conditioning on the time of the second decision epoch and the state at that epoch. If $\pi,\pi^{\dagger} \in \mathcal{D}$ and $\pi^{"} \in P$ are such that $\pi^{"}$ uses π^{\dagger} at the first decision epoch and π thereafter, then

$$\varphi_{\pi^{n}}(s) = \sum_{s' \in S} q_{\pi'}(s,s') \cdot \varphi_{\pi}(s') ,$$

$$w_{\pi^{ii}}(s) = c_{\pi^{i}}(s) - \sum_{s' \in S} \tau_{\pi^{i}}(s,s') \cdot \varphi_{\pi}(s') + \sum_{s' \in S} q_{\pi^{i}}(s,s') \cdot w_{\pi}(s') ,$$

for $s \in S$. If $\varphi_{\pi''}(s) \leq \varphi_{\pi}(s)$ and $w_{\pi''}(s) \leq w_{\pi}(s)$ for $s \in S$, and if, in addition, $\varphi_{\pi''}(s) < \varphi_{\pi}(s)$ or $w_{\pi''}(s) < w_{\pi}(s)$ for some $s \in S$, then π'' is an improvement over π . It can be shown that π'' is also an improvement over π in that case (see Denardo and Fox (1968)). This motivates the following definitions.

A policy π is called <u>unimprovable</u> if

 $\varphi_{\pi}(s) \leq \sum_{s^{i} \in S} q_{\pi^{i}}(s,s^{i}) \cdot \varphi_{\pi}(s^{i}) ,$

$$\mathbf{w}_{\pi}(\mathbf{s}) \leq \mathbf{c}_{\pi^{1}}(\mathbf{s}) - \sum_{\mathbf{s}^{1} \in \mathbf{S}} \tau_{\pi^{1}}(\mathbf{s}, \mathbf{s}^{1}) \cdot \boldsymbol{\varphi}_{\pi}(\mathbf{s}^{1}) + \sum_{\mathbf{s}^{1} \in \mathbf{S}} \mathbf{q}_{\pi^{1}}(\mathbf{s}, \mathbf{s}^{1}) \cdot \mathbf{w}_{\pi}(\mathbf{s}^{1}) ,$$

for $s \in S$ and $\pi^i \in \mathcal{D}$, assuming that all of the expressions above are well-defined and finite. A policy π is <u>strictly</u> unimprovable if it is unimprovable and if, in addition, equalities in the above expression are achieved simultaneously only when $\pi^i = \pi$. If the state space is finite, then an unimprovable policy is average optimal (see Denardo and Fox (1968)). If the state space is not finite, an unimprovable policy is not necessarily average optimal any more (see Hordijk (1974)). Thus, some additional conditions must be satisfied in order to be guaranteed that an unimprovable policy is optimal. Such conditions are given in the next sections.

2. The Case of Constant Optimal Expected Average Cost.

For many semi-Markov decision processes, the optimal long-run expected average cost is constant (i.e., independent of the start-state). In particular, if any state can be reached from each state (by using an appropriate policy) such that the expected cost up to the time the state is reached is well defined and finite, then the optimal long-run expected average cost must be constant. For in this case, the long-run expected average cost, given any start-state s and policy π , can be obtained for any other start-state by using a policy whose actions coincide with those of π at states which are reached from s with a non-zero probability under π , and otherwise are such that the expected cost up to the time when s is reached is finite.

For each $\pi \in \mathcal{D}$, let x_{π} be the mapping from $S \times S$ into R_{+} such that

$$x_{\pi}(s,s') = \lim_{t \to \infty} E_{\pi,s} \{\sum_{n \in \mathbb{N}_{+}} \delta(s_{n},s')\},$$

for s,s' ε S. Here, δ is the Kroenecker delta function, given by

$$\delta(s,s^{\dagger}) = \begin{cases} l & \text{if } s = s^{\dagger} , \\ 0 & \text{if } s \neq s^{\dagger} . \end{cases}$$

The fact that x_{π} exists (although possibly infinite valued) follows from renewal theory (see Smith (1955)). We assume that the expected time until the second decision epoch, given any start-state and action at the first decision epoch, non-zero. This implies that x_{π} is always finite valued.

TAN ANY

Lemma 1: For each $\pi \in \mathfrak{D}$,

$$x_{\pi}(s,s') = \sum_{s'' \in S} x_{\pi}(s,s'') \cdot q_{\pi}(s'',s') ,$$

for $s, s^{t} \in S$.

<u>Proof</u>: For each $\alpha > 0$, $\pi \in \mathfrak{N}$, let $x_{\pi,\alpha}$ be the mapping from $S \times S$ into R_{+} such that

$$x_{\pi,\alpha}(s,s^{*}) = E_{\pi,s}\{\sum_{n \in \mathbb{N}} e^{-\alpha t_n} \cdot \delta(s_n,s^{*})\},$$

for s,s' \in S. Since x_{π} exists,

$$x_{\pi}(s,s^{\dagger}) = \lim_{\alpha \to 0} \alpha \cdot x_{\pi,\alpha}(s,s^{\dagger}) ,$$

for $s, s' \in S$. Now

$$\mathbf{x}_{\pi,\alpha}(s,s') = \sum_{s'' \in S} \mathbf{x}_{\pi,\alpha}(s,s'') \cdot \mathbf{q}_{\pi,\alpha}(s'',s') + (s,s')$$

for $s, s' \in S$, where

$$q_{\pi,\alpha}(s,s') = \int_{t \in \mathbb{R}_+} e^{-\alpha t} dq(s,a_{\pi}(s),s',t)$$
.

This implies that

$$\lim_{\alpha \to 0} \alpha x_{\pi,\alpha}(s,s^{*}) = \lim_{\alpha \to 0} \sum_{s'' \in S} \alpha x_{\pi,\alpha}(s,s'') \cdot q_{\pi,\alpha}(s'',s') ,$$

$$x_{\pi}(s,s^{t}) = \sum_{s^{t} \in S} x_{\pi}(s,s^{t}) \cdot q_{\pi}(s^{t},s^{t}), \text{ for } s,s^{t} \in S$$

Lemma 2: Let ϵ (> 0) be as in Assumption 1. Then, for each $\pi \in \hat{\mathcal{D}}$,

$$\mathbb{E}_{\pi,s}\{\sum_{n\in\mathbb{N}_{t}}\delta(s^{*},s_{n})\}\leq\frac{t}{\epsilon}\cdot\frac{x_{\pi}(s,s^{*})}{x_{\pi}(s,s)}, \text{ for } t\in\mathbb{R}_{+},$$

for states s and s' which are positive recurrent under π .

<u>Proof</u>: Let π be a policy in \mathcal{D} , and let s and s' be positive recurrent states under π . By Lemma 1,

$$x_{\pi}(s,s^{\dagger}) = \sum_{s'' \in S} x_{\pi}(s,s'') \cdot E_{\pi,s}\{\delta(s'',s_2)\} .$$

Using Lemma 1 repeatedly, we obtain

$$x_{\pi}(s,s') = \sum_{s'' \in S} x_{\pi}(s,s'') \cdot E_{\pi,s}\{\delta(s',s_n)\}, \text{ for } n \in \mathbb{N}.$$

Therefore

$$\mathbb{E}_{\pi,s}\{\delta(s^{*},s_{n})\} \leq \frac{x_{\pi}(s,s^{*})}{x_{\pi}(s,s)}, \text{ for } n \in \mathbb{N}.$$

Now

$$\begin{split} \mathbf{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \delta(s^{*}, s_{n}) \} &= \sum_{n \in \mathbb{N}} \mathbf{E}_{\pi,s} \{ \delta(s^{*}, s_{n}) \cdot \mathbb{P}_{\pi,s} \{ \mathbf{t}_{n} \leq \mathbf{t} | s_{n} \} \} \\ &\leq \sum_{n \leq \frac{t}{\epsilon}} \mathbf{E}_{\pi,s} \{ \delta(s^{*}, s_{n}) \} \end{split}$$

by Assumption 1. The lemma now follows by combining the two last results.

20

or

Lemma 3: If π^* is an unimprovable policy such that $\varphi_*(s)$ is constant and

$$\sum_{s^{\dagger} \in S} x_{\pi}(s,s^{\dagger}) \cdot |_{W_{\pi^{*}}(s^{\dagger})}| < \infty ,$$

for s ϵ S and $\pi \epsilon \mathfrak{D}$, then

$$\varphi_{\pi^{*}(s)} \cdot \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot v_{\pi}(s^{*}) \leq \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot c_{\pi}(s^{*}) ,$$

for seS and $\pi \in \mathfrak{D}$.

いたろうないとうというないないないないで

<u>Proof</u>: Let φ be the constant such that $\varphi = \varphi_{\star}(s)$, for $s \in S$. Since π^{\star} is unimprovable,

$$c_{\pi}(s^{*}) \geq w_{\pi^{*}}(s) + \phi \cdot v_{\pi}(s^{*}) - \sum_{s^{''} \in S} q_{\pi}(s^{*},s^{''}) \cdot w_{\pi^{*}}(s^{''}) ,$$

for s' \in S and $\pi \in \mathcal{D}$. Multiplying both sides by $x_{\pi}(s,s')$ and summing up over s' \in S yields

$$\sum_{s' \in S} x_{\pi}(s,s') \cdot c_{\pi}(s')$$

$$\geq \sum_{s' \in S} x_{\pi}(s,s') \{ w_{\pi}(s') + \varphi \cdot v_{\pi}(s') - \sum_{s'' \in S} q_{\pi}(s',s'') w_{\pi}(s'') \} ,$$

for s \in S, $\pi \in \hat{\mathcal{D}}$. The sums on both sides of the above inequality exists, since

$$\sum_{s' \in S} x_{\pi}(s,s') \{ w_{\pi}(s') + \varphi - v_{\pi}(s') - \sum_{s'' \in S} q_{\pi}(s',s'') w_{\pi}(s'') \}^{-}$$

$$\leq \sum_{s' \in S} x_{\pi}(s,s') w_{\pi}(s')^{-} + \sum_{s',s'' \in S} x_{\pi}(s,s') q_{\pi}(s',s'') w_{\pi}(s'')^{+}$$

$$+ \varphi^{-} \cdot \sum_{s' \in S} x_{\pi}(s,s') \cdot v_{\pi}(s')$$

$$= \sum_{s' \in S} x_{\pi}(s,s') W_{\pi}(s') + \sum_{s'' \in S} x_{\pi}(s,s') W_{\pi}(s') + \psi$$

(using Lemma 1 and Lemma 2)

<∞,

using the assumption of the lemma. Now

$$\begin{split} \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot \{w_{\pi}(s^{*}) + \varphi \cdot v_{\pi}(s^{*}) - \sum_{s^{*} \in S} q_{\pi}(s^{*},s^{*})w_{\pi}(s^{*})\} \\ &\geq \sum_{s^{*} \in S} x_{\pi}(s,s^{*})w_{\pi}(s^{*}) + \varphi \cdot \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot v_{\pi}(s^{*}) \\ &- \sum_{s^{*},s^{*} \in S} x_{\pi}(s,s^{*})q_{\pi}(s^{*},s^{*})w_{\pi}(s^{*}) \\ &= \sum_{s^{*} \in S} x_{\pi}(s,s^{*})w_{\pi}(s^{*}) - \sum_{s^{*} \in S} x_{\pi}(s,s^{*})w_{\pi}(s^{*}) + \varphi \cdot \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot v_{\pi}(s^{*}) \\ &\qquad (\text{using Lemma 1}) \\ &= \varphi \cdot \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot v_{\pi}(s^{*}) , \end{split}$$

for s \in S and $\pi \in \mathfrak{D}$, and the lemma follows.

Q.E.D.

For each $\pi \in \mathfrak{D}$, let $R(\pi)$ denote as before the set of positive recurrent states under π , and let $T(\pi)$ denote the set of the other states. For each $\pi \in \mathfrak{D}$, let y_{π} be the mapping from $S \times S$ into R_{+} such that

$$y_{\eta}(s,s^{\prime}) = \begin{cases} E_{\pi,s} \{\sum_{n \in \mathbb{N}} \delta(s^{\prime},s_{n})\}, & \text{for } s^{\prime} \in \mathbb{T}(\pi), s \in S, \\ 1 & \text{for } s^{\prime} \in \mathbb{R}(\pi), s \in S. \end{cases}$$

In words, $y_{\pi}(s,s^{i})$ is the expected number of times the state of the system is s' before a positive recurrent state is <u>entered</u> from another state, given that the start-state is s and that the policy π is used.

Theorem 4: If π^* is an unimprovable policy such that $\varphi_*(s)$ is constant and

$$\sum_{s^{\dagger} \in S} (y_{\pi}(s,s^{\dagger}) + x_{\pi}(s,s^{\dagger})) \cdot |w_{\pi^{*}}(s)| < \infty ,$$

for s ϵ S and π ϵ $\widehat{\mathfrak{I}}$, then π^{*} is average optimal.

Proof: We first show that

$$\varphi_{\pi}(s) \ge \sum_{s' \in S} x_{\pi}(s,s') \cdot c_{\pi}(s')$$
,

for seS and $\pi \in \mathfrak{D}$.

Contra Marine Contra

For each $\pi \in \mathfrak{J}$, let \widetilde{q}_{π} and \widetilde{c}_{π} be the mappings from S imes S and S into R such that

$$\widetilde{q}_{\pi}(s,s^{*}) = \begin{cases} q_{\pi}(s,s^{*}), & \text{for } s \in T(\pi) ,\\ 0 & , & \text{for } s \in R(\pi) , \end{cases}$$
$$\widetilde{c}_{\pi}(s) = \begin{cases} c_{\pi}(s) - \varphi \cdot v_{\pi}(s), & \text{for } s \in T(\pi) ,\\ W_{\pi}(s) & , & \text{for } s \in R(\pi) , \end{cases}$$

for s' ϵ S. Since π^{*} is unimprovable,

$$\widetilde{c}_{\pi}(s) \ge w_{\pi^*}(s) - \sum_{s^* \in S} \widetilde{q}_{\pi}(s,s^*) \cdot w_{\pi^*}(s^*) ,$$

for seS and $\pi \in \mathfrak{D}$. Now

$$\sum_{s \in S} y_{\pi}(s^{"},s) \{w_{\pi^{*}}(s) - \sum_{s^{'} \in S} \tilde{q}_{\pi}(s,s^{'}) \cdot w_{\pi^{*}}(s^{'})\}^{-}$$

$$\leq \sum_{s \in S} y_{\pi}(s^{"},s) \cdot w_{\pi^{*}}(s)^{-} + \sum_{s \in S} y_{\pi}(s^{"},s) \sum_{s^{'} \in S} \tilde{q}_{\pi}(s,s^{'}) w_{\pi^{*}}(s^{'})^{+}$$

$$= \sum_{s \in S} y_{\pi}(s^{"},s) w_{\pi^{*}}(s)^{-} + \sum_{s^{'} \in S} (y_{\pi}(s^{"},s^{'}) - \delta(s^{"},s)) \cdot w_{\pi^{*}}(s^{'})^{+}$$

$$= \sum_{s \in S} y_{\pi}(s^{"},s) w_{\pi^{*}}(s)^{-} + \sum_{s^{'} \in S} y_{\pi}(s^{"},s) \cdot w_{\pi^{*}}(s^{'})^{+} - w_{\pi^{*}}(s^{'})^{+}$$

$$< \infty ,$$

by the last assumption of the theorem. This implies that

$$\sum_{\mathbf{s}\in\mathbf{S}} \mathbf{y}_{\pi}(\mathbf{s}^{"},\mathbf{s})\cdot\widetilde{\mathbf{c}}_{\pi}(\mathbf{s})^{-} < \infty, \quad \mathbf{s}^{"} \in \mathbf{S}, \quad \pi \in \hat{\boldsymbol{\mathcal{D}}}.$$

Thus

A REAL PROPERTY OF A REAL PROPER

$$\sum_{s \in S} y_{\pi}(s",s) \cdot \tilde{c}_{\pi}(s)$$

is well-defined and greater than minus infinity for s" \in S and TF \in $\mathring{\mathcal{D}}$. Now

:

$$\varphi_{\pi}(s) = \lim_{t \to \infty} E_{\pi,s} \{\sum_{n \in \mathbb{N}_{t}} c(s_{n}, a_{n}, t-t_{n})\}/t$$
$$= \lim_{t \to \infty} E_{\pi,s} \{\sum_{n \in \mathbb{N}_{t}} c(s_{n}, a_{n}, \infty)\}/t$$

(by Assumption 2)

$$= \lim_{t \to \infty} \mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}_t} \sum_{s' \in S} \delta(s',s_n) \cdot c_{\pi}(s') \} / t$$

$$= \lim_{t \to \infty} \mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \sum_{s' \in \mathbb{T}(\pi)} \delta(s',s_{n}) \cdot c_{\pi}(s') \} / t$$
$$+ \lim_{t \to \infty} \mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}} \sum_{s' \in \mathbb{R}(\pi)} \delta(s',s_{n}) \cdot c_{\pi}(s') \} / t$$

$$= \varphi + \lim_{t \to \infty} \sum_{s' \in T(\pi)} E_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \delta(s',s_{n}) \} \cdot \widetilde{c}_{\pi}(s')/t$$

+
$$\lim_{t \to \infty} \sum_{s' \in \mathbb{R}(\pi)} E_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \delta(s',s_{n}) \} \cdot (c_{\pi}(s') - \varphi \cdot \nu_{\pi}(s'))/t,$$

using Assumption 2. The first limit is non-negative, since

$$\mathbf{E}_{\pi,s}\{\sum_{n\in\mathbb{N}_{t}}\delta(s^{\prime},s_{n})\}\leq \mathbf{y}_{\pi}(s,s^{\prime}), \text{ for } s^{\prime}\in\mathbb{S},$$

and since

þ

the state of the second se

$$\sum_{\mathbf{s}^{\dagger} \in \mathbf{S}} y_{\pi}(\mathbf{s}, \mathbf{s}^{\dagger}) \cdot \tilde{c}_{\pi}(\mathbf{s}^{\dagger})^{\sim} < \infty .$$

Therefore

$$\varphi_{\pi}(s) \geq \varphi + \lim_{t \to \infty} \sum_{s' \in \mathbb{R}(\pi)} \mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \delta(s', s_{n}) \} (c_{\pi}(s') - \varphi \cdot v_{\pi}(s')) / t .$$

Using Lebesque's bounded convergence theorem, we obtain

$$\lim_{t \to \infty} \sum_{s' \in \mathbb{R}(\pi)} E_{\pi,s} \{ \sum_{n \in \mathbb{N}_{t}} \delta(s',s_{n}) \} \cdot (c_{\pi}(s') - \varphi \cdot v_{\pi}(s')) / t$$

$$= \sum_{s' \in \mathbb{R}(\pi)} x_{\pi}(s,s') \cdot (c_{\pi}(s') - \varphi \cdot v_{\pi}(s')) ,$$

since

$$\begin{split} &\lim_{t \to \infty} \mathbb{E}_{\pi,s} \{\sum_{n \in \mathbb{N}_{t}} \delta(s^{*}, s_{n})\}/t = \tau_{\pi}(s, s^{*}), \text{ for } s^{*} \in S , \\ &\mathbb{E}_{\pi,s} \{\sum_{n \in \mathbb{N}_{t}} \delta(s^{*}, s_{n})\}/t \leq \frac{1}{\epsilon} \cdot \frac{x_{\pi}(s^{*}, s^{*})}{x_{\pi}(s^{*}, s^{*})}, \text{ for } s^{*} \in S, s^{*} \in \mathbb{R}(\pi) , \end{split}$$

and

$$\sum_{\mathbf{s}^{\dagger} \in \mathbf{S}} x_{\pi}(\mathbf{s}^{\dagger}, \mathbf{s}^{\dagger}) (\mathbf{c}_{\pi}(\mathbf{s}^{\dagger}) - \boldsymbol{\varphi} \cdot \boldsymbol{v}_{\pi}(\mathbf{s}^{\dagger}))^{-} < \infty .$$

Thus

$$\varphi_{\pi}(s) \geq \varphi + \sum_{s' \in S} x_{\pi}(s,s')(c_{\pi}(s') - \varphi \cdot v_{\pi}(s')) .$$

Using Lemma 3, we obtain

$$\varphi_{\pi}(s) \geq \varphi$$

Q.E.D.

<u>Corollary 5</u>: Suppose that, for each $s \in S$ and $\pi \in \mathcal{D}$, the expected number of decision epochs occurring before reaching a state in $R(\pi)$ is finite. Then, if π^* is an unimprovable policy such that $\varphi_{\pi}(s)$ is constant and, in addition, $w_{\pi}(s)$ is bounded, then π^* is average optimal.

<u>Proof</u>: In view of the theorem and the fact that $w_{\star}(s)$ is bounded, we only need to show that

$$\sum_{s^* \in S} y_{\pi}(s,s^*) < \infty ,$$

for $s \in S$ and $\pi \in \mathfrak{D}$. But this follows from the first assumption of the corollary, which completes the proof.

Theorem 6: If π^* is a strictly unimprovable policy such that $\varphi_{\pi}(s)$ is constant and, in addition,

$$\sum_{s^{\dagger} \in S} (y_{\pi}(s,s^{\dagger}) + x_{\pi}(s,s^{\dagger})) \cdot |w_{\pi^{*}}(s^{\dagger})| < \infty ,$$

for s ϵ S and $\pi \epsilon \hat{\mathfrak{I}}$, then π^* is undiscounted optimal.

<u>Iroof</u>: Let π be any average optimal stationary, deterministic policy. Following the proof of Lemma 3 and Theorem 4, one can easily see that $a_{\pi}(s) \neq a_{\pi^{*}}(s)$ imply that $\phi_{\pi}(s) > \phi_{\pi^{*}}(s)$ for $s \in R(\pi)$, since π^{*} is strictly unimprovable. This implies that $a_{\pi}(s) = a_{\pi^{*}}(s)$ for $s \in R(\pi)$. From the proof of Theorem 4,

$$c_{\pi}(s) \geq w_{\pi}(s) - \widetilde{q}_{\pi}(s,s')w_{\pi}(s')$$
,

for s \in S. This implies that

$$\sum_{s \in S} y_{\pi}(s",s) \widetilde{c}_{\pi}(s) \geq \sum_{s \in S} y_{\pi}(s",s) \{ w_{\pi}(s) - \sum_{s' \in S} \widetilde{q}_{\pi}(s,s') w_{\pi}(s') \} .$$

It was shown in the proof of Theorem 4 that these sums are well-defined. Now

$$\begin{split} \sum_{s \in S} y_{\pi}(s^{"}, s) \{ w_{\pi^{*}}(s) - \sum_{s^{*} \in S} \widetilde{q}_{\pi}(s, s^{*}) w_{\pi^{*}}(s^{*}) \} \\ &= \sum_{s \in S} y_{\pi}(s^{"}, s) w_{\pi^{*}}(s) - \sum_{s, s^{*} \in S} y_{\pi}(s^{"}, s) \widetilde{q}_{\pi}(s, s^{*}) w_{\pi^{*}}(s^{*}) \\ &= \sum_{s \in S} y_{\pi}(s^{"}, s) w_{\pi^{*}}(s) - \sum_{s^{*} \in S} (y_{\pi}(s^{"}, s^{*}) - \delta(s^{"}, s^{*})) w_{\pi^{*}}(s^{*}) \\ &= w_{\pi^{*}}(s^{"}) , \end{split}$$

for $s'' \in S$. Hence

$$w_{\pi^*}(s^{"}) \leq \sum_{s \in S} y_{\pi}(s^{"},s) \cdot \widetilde{c}_{\pi}(s)$$
,

for s" \in S and $\pi \in \mathfrak{O}$. It is easy to check that

$$w_{\pi}(s^{"}) = \sum_{s \in S} y_{\pi}(s^{"},s) \cdot \widetilde{c}_{\pi}(s)$$
,

for $s^{"} \in S$, so

$$w_{\pi^*}(s^{"}) \leq w_{\pi}(s^{"})$$
 ,

for $s'' \in S$.

A Strangent and the second second

and the state of the second second second

Q.E.D.

<u>Corollary 7</u>: Suppose that for each $s \in S$ and $\pi \in \mathfrak{Y}$, the expected number of decision epochs occurring before reaching a state in $\mathbb{R}(\pi)$ is finite. Then, if π^* is a strictly unimprovable r,licy such that $\phi_{\pi^*}(s)$ is constant and, in addition, $w_{\pi^*}(s)$ is bounded, then π^* is undiscounted optimal.

<u>Proof</u>: The proof proceeds just as in the proof of Corollary 5, and so will not be repeated here.

3. The Case of Non-Constant Optimal Expected Average Cost.

The case when the optimal long-run expected scerage cost varies with the start-state now will be considered. The notation is the same as in Section 2.

Lemma 8: If π^* is a policy such that

$$\begin{split} \varphi_{\pi^{\star}}(s) &\leq \sum_{s^{*} \in S} q_{\pi}(s,s^{*}) \cdot \varphi_{\pi^{\star}}(s^{*}) ,\\ \sum_{s^{*} \in S} x_{\pi}(s,s^{*}) \cdot |\varphi_{\pi^{\star}}(s^{*})| < \infty , \end{split}$$

for $s \in S$ and $\pi \in \mathfrak{D}$, then $\varphi_{*}(s)$ is constant in each positive π recurrent class of states under each policy $\pi \in \mathfrak{D}$.

<u>Proof</u>: Let π be a policy in \mathfrak{O} , and let s be a state in $\mathbb{R}(\pi)$. Using Lemma 1 repeatedly, we obtain

$$\mathbf{x}_{\pi}(s,s") = \sum_{s' \in S} \mathbf{x}_{\pi}(s,s') \cdot \mathbf{E}_{\pi,s}, \{\delta(s",s_n)\},$$

for $n \in \mathbb{N}$ and $s'' \in S$. This implies that

$$\mathbb{E}_{\pi,s} \{\delta(s'',s_n)\} \leq \frac{x_{\pi}(s,s'')}{x_{\pi}(s,s)} ,$$

for $n \in \mathbb{N}$, and s" \in S, since $x_{\pi}(s,s) > 0$. Now

$$\sum_{\mathbf{s}''\in\mathbf{S}}\frac{\mathbf{x}_{\pi}(\mathbf{s},\mathbf{s}'')}{\mathbf{x}_{\pi}(\mathbf{s},\mathbf{s})}\cdot |\boldsymbol{\varphi}_{\pi}(\mathbf{s}'')| < \infty ,$$

11 17

because of the second assumption of the lemma. Using Lebesque's bounded convergence theorem, we obtain

$$\lim_{n \to \infty} \sum_{s'' \in S} E_{\pi,s} \{\delta(s'',s_n)\} \cdot \varphi_{\pi'}(s'')$$
$$= \sum_{s'' \in S} x_{\pi}(s,s'') \cdot \varphi_{\pi'}(s'') ,$$

or equivalently,

$$\lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ \varphi_{\pi}(s_n) \} = \sum_{s'' \in S} x_{\pi}(s,s'') \cdot \varphi_{\pi}(s'') .$$

Let $d^{}_{\overline{\mathcal{T}}}$ be the mapping from S into R such that

$$d_{\pi}(s^{"}) = \sum_{s^{'} \in S} q_{\pi}(s^{"},s^{'}) \cdot \varphi_{*}(s^{'}) - \varphi_{*}(s^{"}) ,$$

for s" \in S. d $_{\pi}$ is well-defined by the first assumption of the lemma. It can easily be shown by induction on n that

$$E_{\pi,s}"\{\sum_{i < n} d_{\pi}(s_{i})\} = E_{\pi,s}"\{\phi_{\pi}(s_{n})\} - \phi_{\pi}(s''),$$

for $n \in N$ and s" \in S. Inserting s for s" in this expression and taking the limits as n goes to infinity, we obtain

$$\lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ \sum_{i \leq n} d_{\pi}(s_i) \} = \lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ \varphi_{\pi}(s_n) \} - \varphi_{\pi}(s) \}$$
$$= \sum_{s' \in S} x_{\pi}(s,s') \cdot \varphi_{\pi^*}(s') - \varphi_{\pi^*}(s) < \infty .$$

The first assumption of the lemma implies that $d_{\pi}(s^{*}) \geq 0$, for $s^{*} \in S$ and $\pi \in \mathfrak{D}$. Using this fact together with

$$\lim_{n \to \infty} E_{\pi,s} \{ \sum_{i < n} d_{\pi}(s_i) \} < \infty ,$$

we obtain $d_{\pi}(s) = 0$. But $s \in R(\pi)$ was chosen arbitrarily, so $d_{\pi}(s) = 0$ for $s \in R(\pi)$. This implies that

$$\varphi_{\pi^{*}}(s) = \sum_{s^{\dagger} \in \mathbb{R}(\pi)} q_{\pi}(s,s^{\dagger}) \cdot \varphi_{\pi^{*}}(s) ,$$

for $s \in R(\pi)$. This, in turn, implies that

$$\varphi_{\pi}^{*}(s) = \lim_{n \to \infty} \mathbb{E}_{\pi,s} \{\varphi_{\pi}^{*}(s_{n})\}$$
$$= \sum_{s' \in S} x_{\pi}^{(s,s')} \varphi_{\pi}^{*}(s')$$

for $s \in R(\pi)$. Now, $x_{\pi}(s,s^{\dagger}) = x_{\pi}(s^{"},s^{\dagger})$, for s and s", if they belong to the same positive recurrent class under π . Thus,

$$\varphi_{\pi}(s) = \sum_{s^{\dagger} \in S} x_{\pi}(s,s^{\dagger}) \cdot \varphi_{\pi}(s^{\dagger}) = \varphi_{\pi}(s^{\prime\prime}) ,$$

for s,s" in the same positive recurrent class under π .

Q.E.D.

For each $\pi \in \hat{\mathcal{D}}$, let $I(\pi)$ be the set of positive recurrent classes, and for each $s \in S$ and $z \in I(\pi)$, let $\mathfrak{F}_{\pi}(s,z)$ be the probability that class z is entered, given start-state s and policy π . <u>Lemma 9</u>: If π^* is an unimprovable policy such that the conditions of the previous lemma hold, and, in addition,

$$\lim_{n \to \infty} \inf \sum_{s' \in T(\pi)} E_{\pi,s} \{\delta(s',s_n)\} \cdot \varphi_{\pi'}(s') \leq 0,$$

for s \in S and $\pi \in \hat{\mathcal{D}}$, then

$$\varphi_{\pi^{\star}}(s) \leq \sum_{z \in I(\pi)} p_{\pi}(s,z) \cdot \varphi_{z}$$
,

for $s \in S$ and $\pi \in \mathfrak{D}$. Here, φ_z is the long-run expected average cost under π^* , given that the start-state is in the class z.

<u>Proof</u>: Let π be any policy in \mathfrak{D} , and let S_z be the set of states belonging to class z for each $z \in I(\pi)$. As in the proof of Lemma 8,

$$\varphi_{\overset{*}{\pi}}(s) \leq \lim_{n \to \infty} \inf E_{\pi,s} \{\varphi_{\overset{*}{\pi}}(s_n)\}$$

 $= \lim_{n \to \infty} \sum_{s^{i} \in \mathbb{R}(\pi)} \mathbb{E}_{\pi,s} \{\delta(s^{i}, s_{n})\} \cdot \varphi_{\pi}(s^{i})$

+
$$\lim_{n \to \infty} \inf \sum_{s' \in \mathbf{T}(\pi)} \mathbb{E}_{\pi,s} \{\delta(s',s_n)\} \cdot \varphi_{\pi^*}(s')$$

$$\leq \lim_{n \to \infty} \sum_{s' \in \mathbf{R}(\pi)} \mathbb{E}_{\pi,s} \{\delta(s',s_n)\} \cdot \varphi_{\pi^*}(s'),$$

 π

for s \in S. The last limit exists and is finite. By Lemma 2,

$$\mathbb{E}_{\pi,s}\{\delta(s^{*},s_{n})\} \leq \mathbb{P}_{\pi}(s,z) \leq \frac{1}{\epsilon} \frac{x_{\pi}(s^{"},s^{*})}{x_{\pi}(s^{"},s^{"})} \text{ for some } s^{"} \in \mathbb{R}(\pi), \epsilon > 0,$$

for s' \in S_z, s \in S and z \in I(π). Now

$$\sum_{z \in I(\pi)} p_{\pi}(s,z) \cdot |\phi_{z}| \leq \frac{1}{\epsilon} \sum_{s'' \in S} \frac{x_{\pi}(s'',s'')}{x_{\pi}(s'',s'')} \cdot |\phi_{*}(s'')| < \infty$$

for s \in S. Therefore, by Lebesgue's bounded convergence theorem,

$$\lim_{n \to \infty} \sum_{s^* \in \mathbb{R}(\pi)} \mathbb{E}_{\pi,s} \{\delta(s^*, s_n)\} \cdot \varphi_{\pi^*}(s^*)$$
$$= \sum_{z \in \mathbb{I}(\pi)} p_{\pi}(s, z) \cdot \varphi_{z^*},$$

for s \in S. We conclude that

$$\varphi_{\pi^*}(s) \leq \sum_{z \in I(\pi)} p_{\pi}(s,z) \cdot \varphi_z$$
,

for $s \in S$.

1 H. L. M.

「「「「「「

Q.E.D.

Let y_{π} be defined as in Section 2.

Theorem 10: If π^* is an unimprovable policy such that

$$\sum_{\substack{s^{\dagger} \in S \\ n \to \infty}} x_{\pi}(s,s^{\dagger}) \cdot |\phi_{\pi^{*}}(s^{\dagger})| < \infty ,$$

$$\sum_{\substack{s^{\dagger} \in S \\ n \to \infty}} \{x_{\pi}(s,s^{\dagger}) + y_{\pi}(s,s^{\dagger})\} \cdot |w_{\pi^{*}}(s^{\dagger})| < \infty ,$$

$$\lim_{\substack{n \to \infty}} \inf \sum_{\substack{s^{\dagger} \in T(\pi) \\ s^{\dagger} \in T(\pi)}} E_{\pi,s}\{\delta(s^{\dagger},s_{n})\} \cdot \phi_{\pi^{*}}(s^{\dagger}) ,$$

for s \in S and $\pi \in \mathfrak{G}$, then π^* is average optimal.

<u>Proof</u>: Let π be any policy in \mathfrak{O} . By Lemma 8, $\varphi_{*}(s)$ is constant π for $s \in S_{z}$ ($z \in I(\pi)$). Therefore, by Theorem 4, $\varphi_{\pi}^{*}(s) \leq \varphi_{\pi}(s)$, for $s \in S_{z}$ ($z \in I(\pi)$). Using Lemma 9 together with this, we obtain $\varphi_{\pi}^{*}(s) \leq \varphi_{\pi}(s)$ for $s \in S$. The costs incurred until a positive recurrent state is reached do not contribute anything to the average cost, since it can be shown (as in the proof of Theorem 4) that the expected cost until a positive recurrent state is reached is finite. Thus, π^{*} must be average optimal.

Q.E.D.

<u>Corollary 11</u>: Suppose that, for each $s \in S$ and $\pi \in \mathcal{D}$, the expected number of decision epochs recurring before reaching a state in $\mathbb{R}(\pi)$ is finite.

If π^* is an unimprovable policy such that $\varphi_*(s)$ and $w_*(s)$ are bounded, then π^* is average optimal.

Proof: We only need to show that

$$\sum_{s' \in S} x_{\pi}(s,s') < \infty ,$$

$$\sum_{s'\in S} y_{\pi}(s,s') < \infty ,$$

for each $s \in S$. The first sum is finite by an assumption made in Section 1, the second sum is finite by Lemma 2, and the third sum is finite by the first assumption of the corollary. Thus, the corollary follows. Q.E.D.

<u>Theorem 12</u>: If π^* is a strictly unimprovable policy such that the conditions of Theorem 10 are satisfied, then π^* is undiscounted optimal. <u>Proof</u>: The proof proceeds just as in the proof of Theorem 6, and so will not be repeated here.

<u>Corollary 13</u>: If π^* is a strictly unimprovable policy such that the conditions of Corollary 11 are satisfied, then π^* is undiscounted optimal.

See the proof of Corollary 11.
CHAPTER 3

SEMI-MARKOV DECISION PROCESSES WITH DISCOUNTING

In this chapter the optimization problem arising when the costs are discounted is investigated. From an economic viewpoint, this problem is somewhat more interesting than the problem without discounting. It has been studied by a number of investigators who have made various assumptions about the state and action spaces, the motion of the system and the costs (see Section 2 in Chapter 1). Here, the assumptions made by other authors are weakened, and more general results are obtained.

In Section 1, there is a formal statement of the problem to be considered. It also contains some preliminary results. In Section 2, some useful operators are introduced. In Section 3, the optimality equation is proven. In Section 4, there are some existence theorems. In Section 5, policy improvement is considered. In Section 6, necessary and sufficient conditions for optimality are presented. Finally, in Section 7, there is an analysis using the contraction properties of a certain operator. An alternative set of necessary and sufficient conditions for optimality are obtained.

1. Problem Formulation.

As before, let S be the state space, $\{A_s\}_{s\in S}$ be the set of action spaces, q be the law of motion, and c be the cost function of the SMDP. For each n in N, let s_n , a_n and t_n denote the state of the system, the action and the time of the nth decision epoch, respectively. The first decision epoch is taken to occur at time zero, so

 $t_1 = 0$. Also, let \mathcal{P} , \mathcal{S} and \mathcal{D} denote the set of all policies, the set of stationary policies and the set of deterministic stationary policies, respectively. Let $A = \bigcup_{s \in S} A_s$.

Let α be a given positive interest rate, and let c_{α} be the mapping from S \times A into R such that

$$c_{\alpha}(s,a) = \int_{t \in \mathbb{R}^+} e^{-\alpha t} dc(s,a,t)$$

for $e \in A_s$ for $s \in S$. In other words, $c_{\alpha}(s,a)$ is the expected discounted cost incurred until the second decision epoch, given that the start-state is s and that the first action is a. Naturally, it is assumed that c_{α} exists.

For each π in \mathcal{O} , let v_{π}^+ , v_{π}^- and v_{π} be the three functions from S into $R_+ \cup \{\infty\}$, $R_+ \cup \{\infty\}$ and $R \cup \{\infty\}$, respectively, such

$$v_{\pi}^{+}(s) = E_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_{n}} \cdot c_{\alpha}(s_{n},a_{n})^{+} \} ,$$

$$v_{\pi}^{-}(s) = E_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_{n}} \cdot c_{\alpha}(s_{n},a_{n})^{-} \} ,$$

$$v_{\pi}(s) = v_{\pi}^{+}(s) - v_{\pi}^{-}(s) ,$$

for s in S, where E is the expectation operator and the subscripts π and s indicate that the start-state is s and that the policy π is used. In words, $v_{\pi}(s)$ is the total expected discounted cost, given that the start-state is s and that the policy π is used. v_{π} is the <u>value function</u> of the policy π . Clearly, v_{π}^+ and v_{π}^- are well-defined (possibly infinite-valued). In order that v_{π} be well-defined, the following assumption is made:

Assumption 1: $v_{\pi}(s) < \infty$, for $s \in S$, $\pi \in \mathcal{O}$.

Let v_{α} be the function from S into $R \bigcup \{\infty, -\infty\}$ such that

$$v_{\alpha}(s) = \inf_{\pi \in \mathbb{Q}} v_{\pi}(s)$$
,

for s in S. For purposes which will become clear later, the following assumption is made.

「日本の市法はいい」の

いっかっているい しい

Assumption 2:
$$v_{\alpha}(s) > -\infty$$
, for $s \in S$.

If there can be an infinite number of decision epochs in a finite amount of time, some of the costs may unintentionally be ignored by the definition of v_{π} . In order to eliminate this problem, the following assumption is made:

<u>Assumption 3</u>: $P_{\pi,s}$ { $t_1 \le t$ for $n \in \mathbb{N}$ } = 0, for $t \in \mathbb{R}_+$, $s \in S$, $\pi \in \mathbb{Q}$.

Here, P is the probability operator and the subscripts π and s indicate that the start-state is s and that the policy π is used. For purposes that will become clear later, a fourth assumption is made:

Assumption 4: Given $\epsilon > 0$, there is an m (possibly depending on s) such that

$$\mathbb{E}_{\pi,s}\{\sum_{n \geq m} e^{-\alpha t} c_{\alpha}(s_{n},a_{n})^{-\gamma}\} \leq \epsilon$$

for π in \mathcal{P} .

These assumptions are satisfied trivially if $c_{\alpha}(s,a)$ is non-negative for each s and a. The following theorem gives some weaker conditions under which the assumptions hold. Theorem 1: If c_{α} is uniformly bounded from below and there is a $\beta < 1$ such that

$$\mathbb{E}_{\pi,s}\{e^{-\alpha t}2\} \leq \beta$$
,

for s in S and π in P, then all the assumptions above hold.

Proof: Let β be as in the theorem. For each $n \in \mathbb{N}$,

$$E_{\pi,s} \{e^{-\alpha t_{n+1}}\}$$

$$= E_{\pi,s} \{e^{-\alpha t_n} \cdot E_{\pi,s} \{e^{-\alpha (t_{n+1}-t_n)} | s_n, a_n\}\}$$

$$\leq \beta \cdot E_{\pi,s} \{e^{-\alpha t_n}\}$$

since

Sales Street

$$\mathbb{E}_{\pi,s}\{e^{-\alpha(t_{n+1}-t_n)}|s_n,a_n\} \leq \beta.$$

This implies that

$$\mathbb{E}_{\pi,s}\{e^{-\alpha t_n}\} \leq \beta^{n-1}$$
,

for $n \in \mathbb{N}$. For each m in \mathbb{N}_{r} .

$$E_{\pi,s} \{ \sum_{n \ge in} e^{-\alpha t_n} \} = \sum_{n \ge m} E_{\pi,s} \{ e^{-\alpha t_n} \}$$
$$\leq \sum_{n \ge m} \beta^{n-1} = \beta^{m-1} \cdot (1-\beta)^{-1} .$$

This implies that

$$(1-\beta)^{-1} \ge E_{\pi,s} \{\sum_{n \in \mathbb{N}} e^{-\alpha t_n}\} \ge E_{\pi,s} \{\sum_{n \le m} e^{-\alpha t_n}\}$$
$$\ge m e^{-\alpha t} \cdot P_{\pi,s} \{t_n \le t \text{ for } n \le m\}$$

for $t \in R_{+}$ and $m \in N$. Thus

$$P_{\pi,s}\{t_n \leq t \text{ for } n \leq m\} \leq e^{\alpha t} \cdot (1-\beta)^{-1} \cdot m^{-1}$$

for t $\in \mathbb{R}_+$ and m $\in \mathbb{N}$. Assumption 3 follows by taking the limits as m goes to infinity.

Let M be an upper bound on $-c_{\alpha}(s,a)$. Then

$$E_{\pi,s} \{\sum_{n \ge m} e^{-\alpha t_n} \cdot c_{\alpha}(s_n, a_n)^{-} \}$$

$$\leq M \cdot E_{\pi,s} \{\sum_{n \ge m} e^{-\alpha t_n} \}$$

$$\leq M \cdot \beta^{m-1} \cdot (1-\beta)^{-1}$$

for $m \in N$. This shows that the rest of the assumptions hold.

2. The Operators Q_{π} and T_{π} .

Let B be the set of mappings from S into $\mathbb{R} \bigcup \{\infty\}$. For each $\pi \in \mathbb{P}$, define the operators \mathbb{Q}_{π} and \mathbb{T}_{π} from B into B by

$$(Q_{\pi}v)(s) = E_{\pi,s} \{e^{-\alpha t_2} \cdot v(s_2)\}, \text{ for } s \in S,$$

$$(T_{\pi}v)(s) = E_{\pi,s} \{c_{\alpha}(s_1,a_1) + e^{-\alpha t_2} \cdot v(s_2)\}, \text{ for } s \in S,$$

for v in B. For some v in B, the above expressions may not be well-defined. Those functions will, however, not be used.

Some compact notation will be used. If u and v are functions in B, then $u \leq v$ means that $u(s) \leq v(s)$ for $s \in S$, u + v is the function such that (u+v)(s) = u(s) + v(s) for $s \in S$, and if c is a constant, then cv is the function such that $(cv)(s) = c \cdot v(s)$ for $s \in S$, etc.

Lemma 2: If u and v in B are such that $u \leq v$, then $Q_{\pi} u \leq Q_{\pi} v$ and $T_{\pi} u \leq T_{\pi} v$, provided the expressions are well-defined.

For each $n \in \mathbb{N}$ and $\pi \in \mathbb{P},$ define the operators \mathbb{Q}_π^n and \mathbb{T}_π^n by

$$(\mathbb{Q}_{\pi}^{n}\mathbf{v})(s) = \mathbb{E}_{\pi,s} \{ e^{-\alpha t} \cdot \mathbf{v}(s_{n+1}) \}, \text{ for } s \in S ,$$

$$(\mathbf{T}_{\pi}^{n}\mathbf{v})(s) = \mathbb{E}_{\pi,s} \{ \sum_{i \leq n} e^{-\alpha t} \cdot \mathbf{c}_{\alpha}(s_{i},a_{i}) + e^{-\alpha t} \cdot \mathbf{v}(s_{n+1}) \}, \text{ for } s \in S ,$$

for v in B. Again, these expressions need not always be well-defined for each v in B. If, however, v is the value function of a policy, then the expressions are clearly well-defined.

Let θ be the function (from S into R) which is zero everywhere. Then

$$(\mathbf{T}_{\pi}^{n}\theta)(\mathbf{s}) = \mathbf{E}_{\pi,\mathbf{s}} \{ \sum_{i \leq n} e^{-\alpha t_{i}} \cdot \mathbf{c}_{\alpha}(\mathbf{s}_{i},\mathbf{a}_{i}), \quad \forall \mathbf{s} \in \mathbf{S},$$

and

$$v_{\pi} = \lim_{n \to \infty} T_{\pi}^{n} \theta$$
,

for any π in \mathfrak{P} .

Lemma 3: For each n ε N and $\pi \in \mathtt{T}, \ \mathtt{Q}_{\pi}^n \mathtt{v}_{\alpha}$ and $\mathtt{T}_{\pi}^n \mathtt{v}_{\alpha}$ are well-defined.

<u>Proof</u>: Let $\epsilon \ge 0$ be given, and let π ' be an ϵ -optimal policy. This means that $v_{\alpha} \ge v_{\pi}$, $+ \epsilon \cdot 1$ where 1 is the function from S into {1}. This implies that

$$\tilde{\mathbf{v}_{\alpha}} \leq (\tilde{\mathbf{v}_{\pi^{\dagger}}} + \epsilon \cdot \mathbf{1})^{-1} \leq \tilde{\mathbf{v}_{\pi^{\dagger}}} + \epsilon \cdot \mathbf{1}$$

Therefore

$$Q_{\pi}^{n} v_{\alpha}^{-} \leq Q_{\pi}^{n} (v_{\pi}^{-}, + \epsilon \cdot 1) \leq Q_{\pi}^{n} v_{\pi}^{-}, + \epsilon \cdot 1$$
,

since

 $Q_{\pi}^{n} \mathbf{1} \leq \mathbf{1}$.

This implies that $Q_{\pi}^n v_{\alpha}^-$ is finite-valued, and thus $Q_{\pi}^n v_{\alpha}$ is well-defined. Also

$$\mathbb{T}^{n}_{\pi} v_{\alpha}^{-} \leq \mathbb{T}^{n}_{\pi} v_{\pi}^{-}, + \epsilon \cdot 1$$

since

 $Q_{\pi}^{n} \mathbf{1} \leq \mathbf{1}$.

This implies that $T^n_{\pi'\alpha}$ is finite-valued, and thus $T^n_{\pi'\alpha}$ is well-defined.

Lemma 4: For each $\pi \in \mathbb{P}$,

$$\lim_{n \to \infty} \inf Q_{\pi}^{n} v_{\alpha} \ge 0.$$

<u>Proof</u>: Let $\epsilon > 0$ be given. From the proof of Lemma 3, there is a policy π^{\dagger} such that

$$Q_{\pi}^{n} v_{\alpha}^{-} \leq Q_{\pi}^{n} v_{\pi}^{-}, + \epsilon \cdot 1$$
,

for all π in P. This implies that

$$\lim_{n \to \infty} Q_{\pi}^{n} \overline{v_{\alpha}} \leq \lim_{n \to \infty} Q_{\pi}^{n} \overline{v_{\pi}}^{*} + \epsilon \cdot 1$$
$$= \epsilon \cdot 1 ,$$

by Assumption 4. The lemma follows, since ϵ is arbitrary.

3. The Optimality Equation.

Bellman (1957) introduced the <u>principle of optimality</u> for dynamic programming. He says (p. 83), "An optimal policy has the property that whatever the initial state and initial decisions are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision." Since an optimal policy need not always exist, the principle has a limited potential use. More useful is the <u>optimality equation</u>, given in the theorem below. For a discussion of the principle of optimality and the optimality equation, see Porteus (1975a).

Let $q_{\gamma\gamma}$ be the mapping from $S \times A \times S$ into R such that

$$q_{\alpha}(s,a,s') = \int_{0}^{\infty} e^{-\alpha t} dq(s,a,s',t) ,$$

for $a \in A_s$ for s', $s \in S$.

Theorem 5: For each s in S,

$$v_{\alpha}(s) = \inf \{c_{\alpha}(s,a) + \sum_{s' \in S} q_{\alpha}(s,a,s') \cdot v_{\alpha}(s')\}.$$

<u>Proof</u>: The proof is similar to the one given in Ross (1970, p. 121) for the case when the action spaces are finite. Let π^i be an ϵ optimal policy. This exists for each $\epsilon > 0$, since $v_{\alpha}(s) > -\infty$ for each $s \in S$ by Assumption 2. Then

$$\mathbf{v}_{\alpha} \leq \mathbf{T}_{\pi} \mathbf{v}_{\pi}, \leq \mathbf{T}_{\pi} (\mathbf{v}_{\alpha} + \epsilon \cdot \mathbf{1}) = \mathbf{T}_{\pi} \mathbf{v}_{\alpha} + \mathbf{Q}_{\pi} (\epsilon \cdot \mathbf{1})$$
$$\leq \mathbf{T}_{\pi} \mathbf{v}_{2} + \epsilon \cdot \mathbf{1}$$

for all $\pi \in \mathbb{Q}$. Since ϵ is arbitrary, $v_{\alpha} \leq T_{\pi}v_{\alpha}$ for $\pi \in \mathbb{Q}$. This is equivalent to

$$v_{\alpha}(s) \leq \inf_{a \in A_{s}} \{c_{\alpha}(s,a) + \sum_{s' \in S} q_{\alpha}(s,a,s') \cdot v_{\alpha}(s')\},$$

for s \in S. We now show that this inequality also holds in the opposite direction.

For each $s \in S$,

Now

$$E_{\pi,s} \{\sum_{n > 1} e^{-\alpha(t_n - t_1)} \cdot c_{\alpha}(s_n, a_n) | a_1, s_2, t_2\} \ge v_{\alpha}(s_2)$$
.

To see this, suppose the opposite. Then there must be a^{i} , s^{i} and t^{i} such that

$$E_{\pi,s} \{\sum_{n > 1} e^{-\alpha(t_n - t_1)} \cdot c_{\alpha}(s_n, a_n) | a_1 = a^{\dagger}, s_2 = s^{\dagger}, t_2 = t^{\dagger} \} < v_{\alpha}(s^{\dagger})$$

For each $n \in N$, let h_n denote the history of the process up to the n^{th} decision epoch (including the state at that time). Let π ' be a policy such that for each history h,

$$P_{\pi',s'}\{a_n = a | h_n = h\} = P_{\pi,s}\{a_{n+1} = a | h_{n+1} = (a',s',t',h)\}$$
.

Then

$$v_{\pi^{i}}(s^{i}) = E_{\pi,s} \{\sum_{n > 1} e^{-\alpha(t_{n}-t_{1})} \cdot c_{\alpha}(s_{n},a_{n}) | a_{1} = a^{i}, s_{2} = s^{i}, t_{2} = t^{i} \}$$

< $v_{\alpha}(s^{i})$,

which is a contradiction. Therefore

$$v_{\pi}(s) \geq E_{\pi,s} \{ c_{\alpha}(s_{1},a_{1}) + e^{-\alpha t_{2}} \cdot v_{\alpha}(s_{2}) \}$$
$$= T_{\pi} v_{\alpha}(s) \cdot$$

This implies that

$$\mathbf{v}_{\pi}(s) \geq \inf_{a \in A_s} \{ \mathbf{c}_{\alpha}(s, a) + \sum_{s' \in S} \mathbf{q}_{\alpha}(s, a, s') \cdot \mathbf{v}_{\alpha}(s') \} ,$$

for s ϵ S. But this holds for each π in \mathbb{P} , so

$$v_{\alpha}(s) \ge \inf_{a \in A_{s}} \{c_{\alpha}(s,a) + \sum_{s' \in S} q_{\alpha}(s,a,s') \cdot v_{\alpha}(s')\},$$

for s ϵ S. Combining this with the result above, the theorem follows.

4. On the Existence of Stationary Optimal and Stationary ϵ -Optimal Policies.

In this section the existence of stationary optimal and stationary

e-optimal policies is investigated. It is important to distinguish between stationary optimal policies and optimal stationary policies. While the former policies are truly optimal, the latter ones are only optimal in the class of stationary policies. Conditions are given for optimal stationary policies to be stationary optimal policies.

<u>Theorem 6</u>: If π is a stationary policy such that $v_{\alpha} = T_{\pi}v_{\alpha}$, then π is optimal.

Proof: Since π is stationary, we obtain

$$v_{\alpha} = T_{\pi}^{n} v_{\alpha}$$
,

$$v_{\alpha} = \lim_{n \to \infty} T_{\pi}^{n} v_{\alpha} = \lim_{n \to \infty} \{T_{\pi}^{n} \theta + Q_{\pi}^{n} v_{\alpha}\}$$
$$\geq \lim_{n \to \infty} T_{\pi}^{n} \theta + \lim_{n \to \infty} \inf Q_{\pi}^{n} v_{\alpha}$$

by Lemma 4. Thus, π is optimal.

The second second second second

 $= v_{\pi}$,

<u>Corollary 7</u>: If each A_s is finite, then there is a stationary optimal policy.

<u>Proof</u>: The existence of a policy π as in the theorem is in this case guaranteed by the optimality equation. <u>Corollary 8</u>: If there is an optimal policy, then there is one which is stationary. and the second second

<u>Proof</u>: Let π be an optimal policy. From the proof of the optimality equation, $v_{\pi} \geq T_{\pi} v_{\alpha}$. Since π is optimal, we obtain $T_{\pi} v_{\alpha} \leq v_{\alpha}$. But $v_{\alpha} \leq T_{\pi} v_{\alpha}$ for all $\pi' \in P$, so $v_{\alpha} = T_{\pi} v_{\alpha}$. Let π'' be the stationary policy such that $T_{\pi''} \equiv T_{\pi}$. By the theorem, π'' is optimal. Thus, there is a stationary optimal policy.

Theorem 9: If for each $s, s^{\dagger} \in S$,

$$\mathbb{E}_{\pi,s}\{\sum_{n\in\mathbb{N}} e^{-\alpha t_n} \cdot \delta(s_n, s^{\dagger})\}$$

is uniformly bounded, then an optimal stationary policy is a stationary optimal policy.

<u>Proof</u>: For each $s, s^{i} \in S$, let $M(s, s^{i})$ be an upper bound on

$$\mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_n} \cdot \delta(s_n, s') \} .$$

Let $\epsilon > 0$ be given. Let v be a mapping from S into R₊ such that $v(s^i) > 0$ for $s^i \in S$ and

$$\sum_{s^{\dagger} \in S} M(s,s^{\dagger}) \cdot v(s^{\dagger}) < \infty ,$$

where s is an element of S. Let π be a stationary policy such that

$$T_{\pi}v_{\alpha} \leq v_{\alpha} + \epsilon \cdot v$$
.

sides of this inequality repeatedly, we obtain

$$\mathbf{T}_{\pi}^{n}\mathbf{v}_{\alpha} \leq \mathbf{v}_{\alpha} + \boldsymbol{\epsilon} \sum_{\mathbf{i}} \sum_{\mathbf{i}} \mathbf{q}_{\pi}^{\mathbf{i}}\mathbf{v} ,$$

for $n \in \mathbb{N}$. Letting n go to infinity in the expression above, we obtain

$$\lim_{n \to \infty} \inf \mathbf{T}_{\pi}^{n} \mathbf{v}_{\alpha} \leq \mathbf{v}_{\alpha} + \epsilon \sum_{n \in \mathbb{N}} \mathbf{Q}_{\pi}^{n} \mathbf{v} \cdot$$

Now

$$\begin{split} \lim_{n \to \infty} \inf \mathbb{T}_{\pi}^{n} \mathbb{v}_{\alpha} &= \lim_{n \to \infty} \inf \{\mathbb{T}_{\pi}^{n} \theta + \mathbb{Q}_{\pi}^{n} \mathbb{v}_{\alpha}\} \\ &= \lim_{n \to \infty} \mathbb{T}_{\pi}^{n} \theta + \lim_{n \to \infty} \inf \mathbb{Q}_{\pi}^{n} \mathbb{v}_{\alpha} \\ &\geq \mathbb{v}_{\pi} , \end{split}$$

by Lemma 4. Thus

$$\mathbf{v}_{\pi} \leq \mathbf{v}_{\alpha} + \epsilon \sum_{n \in \mathbb{N}} \mathbf{Q}_{\pi}^{n} \mathbf{v}$$
 ,

and in particular,

$$v_{\pi}(s) \leq v_{\alpha}(s) + \epsilon \sum_{n \in \mathbb{N}} (Q_{\pi}^{n}v)(s)$$

$$\leq v_{\alpha}(s) + \epsilon \sum_{s' \in S} M(s,s') \cdot v(s')$$
.

Let π be an optimal stationary policy. From above,

$$v_{\pi^{\dagger}}(s) \leq v_{\alpha}(s) + \epsilon \sum_{s^{\dagger} \in S} M(s,s^{\dagger}) \cdot v(s^{\dagger})$$
.

But $\varepsilon > 0$ is arbitrary and $\sum\limits_{s^{\, \prime} \in S} M(s,s^{\, \prime})v(s)$ is finite, so

 $v_{\pi^*}(s) \leq v_{\alpha}(s)$. The argument can be repeated for each $s \in S$, so π^* must be optimal.

Theorem 10: If for each $s' \in S$,

$$\mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_n} \cdot \delta(s_n, s') \}$$

is uniformly bounded, then there are stationary ϵ -optimal policies for all $\epsilon > 0$.

Proof: For each $s' \in S$, let M(s') be a bound on

$$\mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_n} \cdot \delta(s_n,s') \} .$$

Following the proof of the previous theorem, we obtain

$$v_{\pi}(s) \leq v_{\alpha}(s) + \epsilon \sum_{s' \in S} M(s') \cdot v(s')$$
,

for some stationary policy π . Since $\epsilon > 0$ is arbitrary and

$$\sum_{s' \in S} M(s') \cdot v(s') < \infty ,$$

the theorem follows directly.

Corollary 11: If there is a $\beta < 1$ such that

$$E_{\pi,s} \{e^{-\alpha t_2}\} \leq \beta$$
,

for $s \in S$ and $\pi \in \mathbb{P}$, then there are stationary ϵ -optimal policies for arbitrarily small ϵ and every optimal stationary policy is a stationary optimal policy. <u>Proof</u>: We only need to show that the conditions of the two previous theorems are satisfied. It is enough to show that

$$\mathbb{E}_{\pi,s}\{\sum_{n\in\mathbb{N}} e^{-\alpha t_n}\} \leq (1-\beta)^{-1}$$

for s ϵ S, $\pi \in \mathbb{P}$. But this follows from the proof of Theorem 1.

5. Policy Improvement

By the optimality equation, an optimal policy, π^{*} , must satisfy the equation

$$\mathbf{v}_{\pi} = \inf_{\pi \in \mathbf{P}} \mathbf{T}_{\pi} \mathbf{v}_{\pi} \cdot \mathbf{v}_{\pi}$$

A policy π^* which satisfies this equation is called <u>unimprovable</u>. An unimprovable policy need not be optimal, as the following example shows.

Example: Consider a discrete time Markov decision process with state space N_0 and action space {0,1}. If the system enters state 0, the process stops. In every other state there are two permissible actions, 0 and 1. If action 0 is taken in state i(i > 0), an immediate cost a^i is incurred and the state 0 is entered. If action 1 is taken in state i(i > 0), no immediate costs are incurred and the state i + 1 is entered. Let β be a given discount factor. If $a\beta > 1$, then the policy which always chooses action 0 is unimprovable, but it is not optimal.

If there is an $s \in S$ such that

$$v_{\pi^*}(s) > \inf_{\pi \in \mathbb{Q}} (\mathbb{T}_{\pi^v}v_{\pi^*})(s)$$
,

then π^i is called improvable.

Theorem 12: If $\pi' \in \Theta$ and $\pi \in \mathscr{S}$ are such that $\mathbb{T}_{\pi} \mathbf{v}_{\pi'} \leq \mathbf{v}_{\pi'}$, then $\mathbf{v}_{\pi} \leq \mathbb{T}_{\pi} \mathbf{v}_{\pi'}$.

<u>Proof</u>: Applying T_{π} on both sides of $T_{\pi}v_{\pi}, \leq v_{\pi}$ repeatedly, we obtain

 $\mathtt{T}_{\pi} \mathtt{v}_{\pi} \mathtt{,} \geq \mathtt{T}_{\pi}^{n} \mathtt{v}_{\pi} \mathtt{,}$,

for $n \in N$. Letting n go to infinity yields

-

$$\begin{split} \mathbf{T}_{\pi}\mathbf{v}_{\pi^{\dagger}} &\geq \lim_{n \to \infty} \inf \mathbf{T}_{\pi}^{n} \mathbf{v}_{\pi^{\dagger}} = \lim_{n \to \infty} \inf \{\mathbf{T}_{\pi}^{n} \theta + \mathbf{Q}_{\pi}^{n} \mathbf{v}_{\pi^{\dagger}}\} \\ &= \lim_{n \to \infty} \mathbf{T}_{\pi}^{n} \theta + \lim_{n \to \infty} \inf \mathbf{Q}_{\pi}^{n} \mathbf{v}_{\pi^{\dagger}} \\ &\geq \mathbf{v}_{\pi^{\dagger}}, \end{split}$$

by Lemma 4. Thus, the theorem is proved.

This theorem may be useful for the development of a policy improvement procedure like that of Howard (1960). The problem is that one has to avoid convergence to a suboptimal solution.

6. Necessary and Sufficient Conditions for Optimality.

In Section 5, it was shown that an unimprovable policy need not always be optimal. Here, necessary and sufficient conditions for a policy to be optimal are presented. If v_{α} is known, then the optimality equation can be used to find out whether a given policy is optimal or not. If v_{α} is not known in advance, the following theorems may be more useful for proving that a given policy is optimal.

<u>Theorem 13</u>: Let S' be the set of s in S for which $v_{\alpha}(s)$ is

finite. Let \mathcal{P}^{i} be any subset of \mathcal{P} such that for each π in \mathcal{P} there is a π^{i} in \mathcal{P} such that $v_{\pi^{i}} \leq v_{\pi}$.

If π^{\star} is an unimprovable policy such that

$$\lim_{n \to \infty} (Q_{\pi}^{n} v_{\pi})(s) = 0 \text{ for } s \in S', \pi \in \mathcal{O}',$$

then π^* is optimal.

<u>Proof</u>: We first prove that $v_{\pi^*} \leq T_{\pi^*}^n v_{\pi^*}$ for $n \in \mathbb{N}$ and $\pi \in \mathbb{P}$. This clearly holds for n = 1. Now

$$(\mathbb{T}_{\pi}^{n} \mathbb{v}_{\pi}^{*})(s) = \mathbb{E}_{\pi,s} \{ \sum_{i \leq n} e^{-\alpha t_{i}} \cdot c_{\alpha}(s_{i}, a_{i}) + e^{-\alpha t_{n+1}} \cdot \mathbb{v}_{\pi}^{*}(s_{n+1}) \}$$
$$= \mathbb{E}_{\pi,s} \{ \sum_{i < n} e^{-\alpha t_{i}} \cdot c_{\alpha}(s_{i}, a_{i})$$
$$+ e^{-\alpha t_{n}} \cdot \mathbb{E}_{\pi,s}(c_{\alpha}(s_{n}, a_{n}) + e^{-\alpha (t_{n+1} - t_{n})} \cdot \mathbb{v}_{\pi}^{*}(s_{n+1}) | t_{n}, s_{n} \} \}$$

Furthermore,

in an Onio at my anter the spin of the state of the

$$E_{\pi,s} \{ c_{\alpha}(s_{n},a_{n}) + e^{-\alpha(t_{n+1}-t_{n})} \cdot v_{\pi}(s_{n+1}) | t_{n},s_{n} \} \ge v_{\pi}(s_{n}) ,$$

since π^* is unimprovable. Therefore

$$(\mathbb{T}^{n} \mathbf{v}_{\pi^{*}})(s) \geq \mathbb{E}_{\pi,s} \{ \begin{array}{c} \sum \\ i \leq n \end{array} e^{-\alpha t} \cdot \mathbf{v}_{\alpha}(s_{i}, a_{i}) + e^{-\alpha t} \cdot \mathbf{v}_{\pi^{*}}(s_{n}) \} \\ \geq \mathbb{E}_{\pi,s} \{ \begin{array}{c} \sum \\ i \leq n-1 \end{array} e^{-\alpha t} \cdot \mathbf{v}_{\alpha}(s_{i}, a_{i}) + e^{-\alpha t} - 1 \cdot \mathbf{v}_{\pi^{*}}(s_{n-1}) \} \\ \vdots \\ \vdots \\ \geq \mathbb{V}_{\pi^{*}}(s) , \end{array} \right)$$

for s ϵ S. This implies that

$$\begin{array}{l} \mathbf{v}_{\pi} \leq \lim_{n \to \infty} \mathbf{T}_{\pi}^{n} \mathbf{v}_{\pi} = \lim_{n \to \infty} \left\{ \mathbf{T}_{\pi}^{n} \theta + \mathbf{Q}_{\pi}^{n} \mathbf{v}_{\pi} \right\} \\ \\ = \lim_{n \to \infty} \mathbf{T}_{\pi}^{n} \theta + \lim_{n \to \infty} \mathbf{Q}_{\pi}^{n} \mathbf{v}_{\pi} \\ \\ = \mathbf{v}_{\pi} , \end{array}$$

by the last condition of the theorem. Thus, π^{\star} is optimal.

Corollary 14: If π^* is an unimprovable policy such that

$$\sum_{s^* \in S} E_{\pi,s} \{ e^{n} \cdot \delta(s^*, s_n) \} \cdot |v_{\pi^*}(s^*)| < \infty, \text{ for } s \in S^*, \pi \in \mathbb{P}^* \},$$

then π^{\star} is optimal.

Proof: It is sufficient to recognize that

$$\sum_{s' \in S} E_{\pi,s} \{ e^{-\alpha t_n} \cdot \delta(s', s_n) \} \cdot |v_{\pi'}(s')|$$
$$= \sum_{n \in \mathbb{N}} (Q_{\pi}^n |v_{\pi'}|)(s)$$
$$\geq \sum_{n \in \mathbb{N}} |(Q_{\pi'}^n v_{\pi'})(s)| ,$$

for $s \in S'$ and $\pi \in \mathbb{Q}^{\circ}$. Thus the conditions of the theorem hold, and the corollary follows.

<u>Corollary 15</u>: If π^* is an unimprovable policy such that its valuefunction is bounded, then π^* is optimal.

<u>Proof</u>: Let M be an upper bound on $|v_{\pi}(s)|$. Then

$$\lim_{n \to \infty} |E_{\pi,s} \{e^{-\alpha t_n} \cdot v_*(s_n)\}|$$
$$= \lim_{n \to \infty} E_{\pi,s} \{e^{-\alpha t_n} \cdot |v_*(s)|\}$$
$$\leq M \cdot \lim_{n \to \infty} E_{\pi,s} \{e^{-\alpha t_n}\} = 0,$$

by Assumption 3. The corollary now follows from the theorem.

<u>Theorem 16</u>: Suppose that there is an optimal policy, π . Then a policy π^* is optimal if and only if

$$\lim_{n \to \infty} (Q_{\pi}^{n} v_{\star})(s) = 0, \text{ for } s \in S^{*}.$$

$$(Q_{\pi}^{n}v_{\pi}^{*})(s) = (Q_{\pi}^{n}v_{\pi})(s) = v_{\pi}(s) - (T_{\pi}^{n}\theta)(s)$$
,

for $s \in S'$, $n \in N$. This implies that

$$\lim_{n \to \infty} (Q_{\pi}^{n} v_{\star})(s) = \lim_{n \to \infty} \{v_{\pi}(s) - (T_{\pi}^{n} \theta)(s)\} = 0,$$

for s \in S'. This completes the proof of the theorem.

7. Norms and Contraction Mar _ings.

It may sometimes be more convenient to work with norms and contraction mappings. Denardo (1967) did this, and developed an elegant analysis. Recently, Lippman (1975) used these concepts.

As before, let 1 be the function from S into R with value 1 everywhere. Let $\|\cdot\|$ be a norm on B such that

(a) $\|\underline{1}\| = 1$, (b) $\|u\| \le \|v\|$ if $0 \le u \le v$.

The sup norm, given by

A STATE

一川南川

$$\|v\| = \sup_{s \in S} |v(s)|$$
,

is such a norm. Lippman (1975) has considered other norms.

A mapping T from B into B is called a contraction mapping if there is a $\beta < 1$ such that

$$||Tv|| \leq \beta ||v||$$

for $v \in B$. Denardo's <u>n-stage</u> contraction condition is as follows. There is an $n \in N$ and a $\beta < 1$ such that

$$\|Q_{rr}^{n}v\| \leq \beta \|v\|$$
 ,

for $v \in B$ and $\pi \in \mathbb{P}$. We weaken the n-stage contraction condition so that it reads as follows. For each $v \ge 0$, there is an n N and a $\beta < 1$ such that

$$\|\mathbf{Q}_{\pi}^{n}\mathbf{v}\| \leq \beta \|\mathbf{v}\| ,$$

for all π in \mathfrak{F} .

Lemma 17: If there is an $n \in \mathbb{N}$ and a $\beta < 1$ such that

$$\mathbb{E}_{\pi,s}\{e^{-\alpha t}n\} \leq \beta$$

for $s \in S$ and $\pi \in \mathcal{P}$, then the sup norm satisfies the n-stage contraction condition. Proof: We have

$$\begin{aligned} |Q_{\pi}^{n}v|| &= \sup_{s \in S} |(Q_{\pi}^{n}v)(s)| &= \sup_{s \in S} |E_{\pi,s} \{e^{-\alpha t}n v(s_{n})\}| \\ &\leq \sup_{s \in S} |E_{\pi,s} \{e^{-\alpha t}n \cdot \sup_{s \in S} |v(s')|\} \\ &= ||v|| \cdot \sup_{s \in S} |E_{\pi,s} \{e^{-\alpha t}n\}, \end{aligned}$$

and the lemma follows.

and the second state of the second second

And the second se

Let $\rho(\cdot, \cdot)$ be a metric on $B \times B$ such that for u,v in B, $\rho(u,v) = ||w||$, where

$$w(s) = \begin{cases} u(s) - v(s), & \text{if } u(s) < \infty \text{ or } v(s) < \infty \text{,} \\ 0, & \text{if } u(s) = v(s) = \infty \text{.} \end{cases}$$

<u>Theorem 18</u>: If $\|\cdot\|$ satisfies the n-stage contraction condition, then a policy π^* is optimal if and only if π^* is unimprovable and $\rho(v_{\pi^*}, v_{\alpha}) < \infty$.

<u>Proof</u>: The only if part of the theorem is trivial. We now prove the if part. Let w be such that

$$w(s) = \begin{cases} v_{\pi}(s) - v_{\alpha}(s), & \text{if } v_{\alpha}(s) < \infty \text{ or } v_{\pi}(s) < \infty, \\ \pi & \pi \\ 0, & \text{if } v_{\pi}(s) = v_{\alpha}(s) = \infty. \end{cases}$$

Let $n \in \mathbb{N}$ and $\beta < 1$ be as in the contraction condition. Let $\epsilon > 0$ be given, and let π be a stationary policy such that

$$\mathbb{T}_{\pi} \mathbf{v}_{\alpha} \leq \mathbf{v}_{\alpha} + (\epsilon/n) \cdot \mathbf{1}$$
.

 π^* is unimprovable, so $v_{\pi^*} \leq T_{\pi^v} v_{\pi^*}$. This implies that $w \leq Q_{\pi^w} + (\epsilon/n) \cdot 1$,

and the second states of the second second

since w $\geq \theta.$ Applying $\mathbf{Q}_{_{\rm T\!\!T}}$ to both sides of this inequality repeatedly yields

$$w \leq Q_{\pi}^{n} w + (\epsilon/n) \sum_{i < n} Q_{\pi}^{i} \mathbf{1}$$
$$\leq Q_{\pi}^{n} w + \epsilon \cdot \mathbf{1},$$

since $Q_{TT} \leq 1$. Taking the norm of the functions on both sides of the inequality yields

$$\begin{split} \|\mathbf{w}\| &\leq \|\mathbf{Q}_{\pi}^{n}\mathbf{w} + \boldsymbol{\varepsilon} \cdot \mathbf{1}\| \\ &\leq \|\mathbf{Q}_{\pi}^{n}\mathbf{w}\| + \boldsymbol{\varepsilon}\| \mathbf{1}\| \\ &\leq \beta \|\mathbf{w}\| + \boldsymbol{\varepsilon} , \end{split}$$

by the contraction condition. But $\epsilon > 0$ is arbitrary, $\beta < 1$ and $||w|| < \infty$. This implies that ||w|| = 0, and thus $v_{\pi^*} = v_{\alpha}$. <u>Corollary 19</u>: If π^* is an unimprovable policy such that $v_{\pi^*}(s) - v_{\alpha}(s)$ is bounded, and if the condition of Lemma 17 is satisfied, then π^* is optimal.

Notice the similarity of this corollary and Corollary 15.

<u>Corollary 20</u>: If $\|\cdot\|$ satisfies the n-stage contraction condition, if v_{α} is bounded below, and if π^* is an unimprovable policy such that $\|v_{\pi^*}\| < \infty$, then π^* is optimal. <u>Proof</u>: Let $-M(M \ge 0)$ be a lower bound on $v_{\alpha}(s)$. Let w be as in the proof of the theorem. Then

14

Same and

$$0 \leq w \leq v_{\pi^*} + M \cdot 1$$
,

so

$$\|w\| \le \|v_{\pi^*} + M \cdot \| \le \|v_{\pi^*}\| + M < \infty$$

CHAPTER 4

OPTIMAL CONTROL OF QUEUEING SYSTEMS

There has been a considerable interest in the control of queueing systems in the last decade. Often the control problems have been formulated in the framework of semi-Markov decision processes. The existence of certain simple and intuitive optimal policies have been proven for many different queueing systems. For a brief (but excellent) survey of the literature in this area, see Gross and Harris (1974, pp. 364-380).

In this chapter, three aspects of the control of queueing systems are considered. In Section 1, the formulation of queueing control problems is discussed. Section 2 elaborates upon two general approaches to the solution of queueing control problems. In Section 3, four different methods for proving the optimality of an unimprovable policy are developed.

1. Formulation of Queueing Control Problems.

The formulation of queueing control problems plays an important role in the solution of these problems. Sometimes, a queueing control problem may be formulated in two different but equivalent ways, where only one is amenable to analysis. Special queueing control problems may have special desirable formulations. But since a general formulation of queueing control problems may yield a better perspective, we shall now briefly describe the various components of a controllable queueing system.

A queueing system consists of an <u>input</u> <u>source</u>, a <u>queue</u> and a <u>service</u> mechanism. The input source generates customers which need certain services

provided by the service mechanism. A customer generated by the input source is said to <u>arrive</u> at the queueing system. The times between two consecutive arrivals are the <u>interarrival times</u>. On arrival, a customer either is given service immediately or is placed in the queue of customers waiting to be served. There may be several <u>customer classes</u>, reflecting the special needs of the customers. The service mechanism may consist of one or several <u>service facilities</u>, each of which has a certain number of <u>servers</u>. When the customers have received their service(s), they <u>leave</u> the system.

-

いいのはいろの

State State

The control of queueing systems can take various forms. Sometimes, the <u>arrival rate</u> may be adjusted dynamically. Other times, the <u>service</u> <u>rate(s)</u> or the <u>number of active servers</u> may be controlled. A third possibility is to control the <u>order</u> in which the customers are given service.

There are various costs that may need to be considered when analyzing queueing systems. For example, there may be a <u>service cost</u> which is incurred each time a customer is served. If the server(s) can be turned on and off, there may be <u>start-up</u> and <u>shut-down costs</u> when the server(s) are turned on and turned off, respectively. There may be an <u>idling cost</u> which is incurred at a positive and constant rate for each server when he is not giving service or performing other useful duties. There may be a <u>customer holding cost</u> which is incurred at a rate which is a function of the number of customers in the system.

There may, of course, be many other types of controls and costs than those which have been mentioned here. But surprisingly many of the queueing control problems which have been considered in the literature fit the above description.

By formulating a queueing control problem as a semi-Markov decision process, the theory for such processes may be used in developing a solution procedure or to prove that a given policy is optimal (or not). The formulation is usually quite straightforward. One only has to define the state of the system and the decision epochs. The state space, the set of action spaces, the law of motion and the cost function of the semi-Markov decision process are then determined by the specification of the queueing system.

The definition of the state of the system is crucial. The state must characterize the queueing system completely at each decision epoch. Since a queueing system consists of an input source, a queue and a service mechanism, one may define the state of the <u>input source</u>, the state of the queue and the state of the <u>service mechanism</u>. The state of the <u>system</u> is then given by these three states. The state space of the system may be defined as the Cartesian product of the state spaces of the input source, the queue and the service mechanism, respectively.

The state space of a queueing system is often countable. If the input source, the queue and the service mechanism all have countable state spaces, then the state of the system is countable.

Consider the state space of the queue. Suppose that there is a countable number of customer classes. If the state of the queue is defined as the vector whose ith component indicates the number of customers in class i (for each $i \in N$), then the state space of the queue is countable. This follows from the fact that there are only a finite number of customers in the queue at any given time.

Consider the state space of the service mechanism. One case is the system which can be controlled by turning serving on or off. For

this case, if there is a countable number of servers, and if the state of the service mechanism is defined as the vector whose ith component indicates whether the ith server is on or off (for each $i \in N$), then the state space of the service mechanism is countable. For a more general case, suppose now that the service rate of each server may be adjusted to a countable number of levels. Also suppose that there are a countable number of servers and that the service rate is only non-zero for a finite number of servers at any given point in time. If the state of the service mechanism is defined as the vector whose ith component indicates the level of the service rate of the ith server, then the state space is still countable.

The definition of the decision epochs is also crucial. As mentioned before, the state of the system must characterize the queueing system completely at each decision epoch. The most natural way to define the decision epochs is by letting them be the epochs when the state of the system changes. If the state of the system (as it happens to be defined) does not characterize the queueing system completely at each of these decision epochs, one can try to eliminate some of the decision epochs.

Sometimes it may be desirable to have the decision epochs equally spaced in time. In this case, the decision epochs are determined by specifying the length of time between two consecutive decision epochs. Magazine (1971) used this approach. Other times, it may be desirable to define the decision epochs such that the times between two consecutive decision epochs are independent and identically distributed random variables. Lippman (1975) used this approach. Both of these ways of defining the decision epochs are motivated by a certain solution method which will be elaborated upon in the next section.

2. Analytical Solution Methods.

A large variety of queueing control problems have been successfully analyzed by a number of investigators. Their successes have to some extent depended on the special features of the problems they considered. But many of the queueing problems also have much in common. Therefore, there is some basis for developing general approaches for solving them. Prabhu and Stidham (1973) attempted to develop a unified view of the different approaches that have been used previously.

If the state and action spaces are finite, then there are wellknown (policy improvement, policy iteration) algorithms for finding an optimal policy. But in the context of queueing systems, one is often more interested in showing that there is an optimal policy of a simple and intuitive form. As a by-product of this, one may perhaps develop especially efficient algorithms for finding an optimal policy. Two general approaches for analyzing queuing systems will now be presented.

The first approach consists of solving the problem for one period (stage) and then extending the results to arbitrarily many periods by an inductive argument. This approach was initially used for solving inventory problems (e.g. by Iglehart (1963)). Because of the similarity between queueing and inventory problems, the approach was later adopted by queueing theoreticians. McGill (1969) used the approach in his analysis of the M/M/c queueing system with controllable servers. A full development of this approach can be found in Porteus (1975b).

This approach has two advantages. First, the one-period problem is usually easier to analyze than the infinite period problem. A successful analysis solves both the finite and infinite horizon problems.

However, this approach of first solving the one-period problem can also have its disadvantages. In fact, for many queueing problems, the one-period problem is rather meaningless. One reason is that the length of the first period may not be nearly the same for different start-states and different actions. Furthermore, many important costs may be neglected in the one-period problem (e.g., switching costs). Nevertheless, the approach is still attractive for many problems.

The second approach consists of restricting one's search for an optimal policy to a small class of stationary policies (hopefully not excluding the optimal policy) and then proving that the policy which is optimal in this class is also optimal among all policies. To prove that a policy believed to be optimal is indeed optimal among all policies, one usually only has to prove that the policy is unimprovable. This approach has been used by, among others, Reed (1974a), (1974b).

This approach has the advantage that it usually only requires the analysis of relatively simple stationary policies. If one can obtain an explicit expression for the value functions of these policies, then it is usually a simple matter to prove when one of these policies is unimprovable (and thus probably optimal). Even if such explicit results cannot be obtained, the approach may still be used with success (e.g., see Orkenyi (1976)).

The disadvantage of the approach lies in the fact that an unimprovable policy need not necessarily be optimal. In the previous chapters, several conditions for an unimprovable policy to be optimal were given. For example, when discounting is used, it was shown that if the value function of the unimprovable policy is bounded, then the policy is optimal.

But queueing control problems are often characterized by giving rise to unbounded value functions. This is often due to the holding costs Leing unbounded. In the next section, it is shown how this problem can be solved.

3. Solutions to the Problem of Unbounded Costs.

We now consider the problem of unbounded costs with discounting, and develop four different methods for proving that an unimprovable policy is optimal. The assumptions of chapter 3 are retained here.

3.1 A Reformulation.

Perhaps the easiest way to solve the problem of unbounded costs is by reformulating the cost structure of the system under consideration in such a way that the costs become bounded. There is, however, no single receipe for doing this. Different problems may require different reformulations. Here, an idea of Bell (1971) is generalized.

For the sake of simplicity, suppose that the expected discounted cost excluding the cost due to holding customers in the system is bounded. Also suppose that there are m customer classes and that a holding cost is incurred at a rate which is a given function, h, of the number of customers present in each customer class. Define the state of the queue as indicated in Section 1.

For each $n \in N$, let t_n denote the time of the n^{th} change in the state of the queue and let y_n denote the state of the queue immediately after the change. Without loss of generality, assume that $t_1 = 0$. For each policy π and state s, let $v_{\pi}^{h}(s)$ denote the expected discounted holding cost, given that the policy π is used and that the

start state is s. Clearly

$$v_{\pi}^{h}(s) = E_{\pi,s} \{ \sum_{n \in \mathbb{N}} \int_{t_{n}}^{t_{n+1}} h(y_{n}) e^{-\alpha t} dt \}$$
$$= \frac{1}{\alpha} h(y_{1}) + E_{\pi,s} \{ \sum_{n \in \mathbb{N}} \frac{1}{\alpha} (h(y_{n+1}) - h(y_{n})) e^{-\alpha t_{n+1}} \},$$

for each seS and $\pi \in \mathbb{P}$.

Now, reformulate the holding cost structure such that at each time $t_n(n > 1)$, the holding cost

$$x_n = \frac{1}{\alpha}(h(y_n) - h(y_{n-1}))$$

is incurred. Formally, we choose to include the cost x_n in the costs incurred in the period from t_{n-1} to $t_n(n > 1)$. For each start-state s and policy π used, the expected discounted holding cost becomes

$$v_{\pi}^{h}(s) - \frac{l}{\alpha}h(y_{l})$$
.

Thus, the problem before the reformulation is equivalent to the problem after the reformulation with regard to optimal policies.

Assume that the number of customers in each customer class only can change by one at a time and that changes in different customer classes cannot occur simultaneously. Let Y denote the state space of the queue, and for each $i(\leq m)$, let ω_i denote the m-vector whose components are all zero except for the ith one which is equal to one. We can now state the following theorem.

Theorem 1: If for each policy π ,



is uniformly bounded, and if there is an $M < \infty$ such that

$$|h(y + u_i) - h(y)| \leq M$$

きまうになったいいれ

for $1 \le i \le m$ and $y \in Y$, then every unimprovable policy is optimal.

<u>Proof</u>: Under the conditions of the theorem, the expected discounted holding cost after the reformulation is bounded. Therefore, any policy which is unimprovable for the problem after the reformulation is optimal for that problem. But the optimal policies are the same for both problems. The unimprovable policies are also the same for both problems. Therefore, we conclude that a policy which is unimprovable for the original problem is also optimal.

Example (The M/G/l queueing system with removable server):

Excluding the policies which turn the server on and off repeatedly at a decision epoch, the expected discounted cost excluding those due to holding customers in the system is bounded. Let λ be the arrival rate of the customers, and let $\omega(<1)$ be the Laplace transform of the service times (with its parameter being equal to the interest rate α).

Let $\{t_n^n\}_{n\in\mathbb{N}}$ be the sequence of times when customers arrive, and let $\{t_n^n\}_{n\in\mathbb{N}}$ be the sequence of times when customers depart. It can easily be shown that for each policy π used and each start-state s,

$$\mathbb{E}_{\pi,s}\{\sum_{n\in\mathbb{N}} e^{-\alpha t_n^*}\} = 1 + \frac{\lambda}{\alpha},$$

and

$$\mathbb{E}_{\pi,s} \{ \sum_{n \in \mathbb{N}} e^{-\alpha t_n''} \} \leq \frac{1}{1-\omega} < \infty .$$

Since $\{t_n\}_{n \in \mathbb{N}}$ is a subsequence of $\{t_n^i\}_{n \in \mathbb{N}} \bigcup \{t_n^u\}_{n \in \mathbb{N}}$, the first condition of the theorem holds.

If the slope of h is bounded (in this case h is a function of one variable), then the second condition of the theorem holds. Thus, if the slope of the holding cost function is bounded, then every unimprovable policy is optimal. This is just the assumption made by Blackburn (1971) when he considered the convex holding cost model.

3.2 Comparison with the Policy which Shuts Down the System.

Assume as before that the customer holding cost is incurred at a rate $h(y_n)$ in each interval (t_n, t_{n+1}) . Also assume that h is such that

$$0 \leq h(x) \leq h(y)$$

for $x \leq y$ and $x \in Y$, $y \in Y$.

Assume that the system can be shut down at any decision epoch and that the shut-down cost is bounded uniformly from above. Let π_0 denote the policy which always shuts the system down (or leaves it off). Assume that when the policy π_0 is used the total number of customers present in each customer class is at a maximum at all times for any given startstate.

<u>Theorem 2</u>: If π^* is an unimprovable policy such that, for each s \in S,

$$v_{\pi^*}(s) \leq v_{\pi_0}(s) < \infty$$
,

then π^{\star} is optimal.

A STATE OF A STATE OF

Proof: By Theorem 13 in Chapter 3, we only need to show that

$$\lim_{n \to \infty} E_{\pi,s} \{ e^{-\alpha t_n} \cdot v_*(s_n) \} = 0$$

for each $s \in S$ and $\pi \in \Theta$. Here $\{t_n\}_{n \in \mathbb{N}}$ is the sequence of the times of the decision epochs.

For each $s \in S$, let R(s) denote the expected discounted shutdown cost when the system is in state s and the policy π_0 is used. For each $\pi \in \Theta$, $s \in S$ and $t \in R$, let $x_{\pi}(s,t)$ denote the discounted holding cost incurred from time t onward (the discounting starting at time 0), given that the start-state is s and that the policy π_0 is used. It follows from our assumptions that

$$x_{\pi}(s,t) \leq x_{\pi_0}(s,t)$$
, for $t \in \mathbb{R}$, $s \in S$, $\pi \in \Theta$.

Now

$$v_{\pi_0}(s) = R(s) + E\{x_{\pi_0}(s,0)\}, \text{ for } s \in S,$$

so

$$E\{x_{\pi_0}(s,0)\} < \infty$$
, for $s \in S$.

For each $\pi \in \mathbb{P}$ and $s \in S$, let $\{t_n(\pi,s)\}_{n \in \mathbb{N}}$ be the sequence of the times of the decision epochs, given that the start-state is s and that the policy π is used.

Choose a $\pi \in \mathbb{P}$, and for each $n \in \mathbb{N}$, let π_n be the policy which follows π until the nth decision epoch and then shuts down the system. Then

$$E_{\pi,s} \{ e^{-\alpha t_n} v_{\pi_0}^{h}(s_n) \} = E\{x_{\pi_n}(s, t_n(\pi, s)) \}$$
$$= E\{ l_{\{t_n(\pi, s) \le t\}} \cdot x_{\pi_n}(s, t_n(\pi, s)) \}$$

$$+ E\{l_{\{t_n(\pi,s) > t\}} \cdot x_{\pi_n}(s,t_n(\pi,s))\}$$

$$\le E\{l_{\{t_n(\pi,s) \le t\}} \cdot x_{\pi_0}(s,0)\}$$

$$+ E\{l_{\{t_n(\pi,s) > t\}} \cdot x_{\pi_0}(s,t_n(\pi,s))\}$$

$$\le E\{l_{\{t_n(\pi,s) \le t\}} \cdot x_{\pi_0}(s,0)\} + E\{x_0(s,t)\}$$

for $n \in \mathbb{N}, \ t \in \mathbb{R}$ and $s \in S.$ Here, we have used the fact that

$$x_{\pi}(s,t) \leq x_{\pi_0}(s,t)$$
, for $t \in \mathbb{R}$, $s \in S$, $\pi \in \mathbb{Q}$

and

$$x_{\pi}(s,t) \leq x_{\pi}(s,t^{i})$$
, for $t \leq t^{i}$, $t,t^{i} \in \mathbb{R}$, $s \in S$, $\pi \in \mathfrak{S}$,

and

$$\mathbf{x}_{\pi}(\mathsf{s},\mathsf{t})\geq \mathsf{0}$$
, for $\mathsf{t}\in\mathsf{R}$, $\mathsf{s}\in\mathsf{S}$, $\pi\in Q$.

By Lebesgue's bounded convergence theorem,

$$\lim_{n \to \infty} \mathbb{E}\{ \mathbb{I}_{\{\mathfrak{t}_n(\pi,s) \leq t\}} \cdot \mathbb{X}_{0}(s,0) \} = 0, \text{ for } s \in S,$$

since

and the second

The Designation of the

$$E\{x_{\pi_0}(s,0)\} < \infty ,$$

and

$$\lim_{n \to \infty} P_{\pi,s} \{ t_n \le t \} = 0 .$$

Thus

$$\lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ e^{-\alpha t} \cdot v^h_{\pi_0}(s_n) \} \le \mathbb{E} \{ x_0(s,t) \}, \text{ for } t \in \mathbb{R}, s \in S .$$

 $\lim_{t \to \infty} E\{x_0(s,t)\} = 0, \text{ for } s \in S,$

since

$$E\{x_{0}(s,0)\} = \lim_{t \to \infty} E_{\pi_{0},s}\{\int_{0}^{t} e^{-\alpha t} \cdot h(y_{t})dt\} < \infty, \text{ for } s \in S,$$

where y_t denotes the state of the queue at time t. Therefore

$$\lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ e^{-\alpha t_n} \cdot \mathbb{V}_{\pi_0}(s_n) \}$$
$$= \lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ e^{-\alpha t_n} \cdot \mathbb{R}(s_n) \} + \lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ e^{-\alpha t_n} \cdot \mathbb{V}_{\pi_0}^h(s_n) \}$$
$$\leq \lim_{n \to \infty} \mathbb{E}_{\pi,s} \{ e^{-\alpha t_n} \cdot \mathbb{R}(s_n) \} ,$$

for $s \in S$. Let M be a finite upper bound on R(s). Then

$$\lim_{n \to \infty} E_{\pi,s} \{ e^{-\alpha t_n} \cdot v_{\pi_0}(s_n) \} \le M \cdot \lim_{n \to \infty} E_{\pi,s} \{ e^{-\alpha t_n} \}$$
$$= 0, \text{ for } s \in S,$$

since

$$\lim_{n \to \infty} \mathbb{P}_{\pi,s} \{ t_n \leq t \} = 0, \text{ for } t \in \mathbb{R}, s \in S.$$

This completes the proof.

Q.E.D.

<u>Example</u> (The M/G/1 queueing system with removable server):

The state of the system is defined as a pair of integers (i,j), where i denotes the number of customers in the system and

But
$$j = \begin{cases} 0, & \text{if the server is off} \\ 1, & \text{if the server is on} \end{cases}$$

It is easy to find that

$$\mathbf{v}_{\pi_{0}}(\mathbf{i},\mathbf{j}) = \begin{cases} \sum_{\mathbf{k}\in\mathbb{N}_{0}}^{\cdot} \left(\frac{\lambda}{\lambda+\alpha}\right)^{\mathbf{k}} \cdot \mathbf{h}(\mathbf{i}+\mathbf{k}), & \text{for } \mathbf{j} = 0, \ \mathbf{i} \in \mathbb{N}_{0} \ , \\\\ \mathbf{R}_{2} + \sum_{\mathbf{k}+\mathbb{N}_{0}}^{\cdot} \left(\frac{\lambda}{\lambda+\alpha}\right)^{\mathbf{k}} \cdot \mathbf{h}(\mathbf{i}+\mathbf{k}), & \text{for } \mathbf{j} = \mathbf{l}, \ \mathbf{i} \in \mathbb{N}_{0} \ . \end{cases}$$

Therefore, if π is an unimprovable policy such that

$$v_{\pi}(i,j) \leq v_{\pi_0}(i,j)$$
, for $j \in \{0,1\}$, $i \in \mathbb{N}_0$,

and if

$$\sum_{\mathbf{i}\in\mathbb{N}_{0}}\left(\frac{\lambda}{\lambda+\alpha}\right)^{\mathbf{i}}\mathbf{h}(\mathbf{i})<\infty,$$

then π is optimal.

3.3 <u>Comparison with the Policy which Minimizes the Expected Discounted</u> Holding Cost.

Suppose that there is a policy which minimizes the expected discounted holding cost, and let π_0 denote such a policy. For each $\pi \in 0^\circ$ and $s \in S$, let $v_{\pi}^{nh}(s)$ denote the expected discounted cost excluding the holding costs, given that the start-state is s and that the policy π is used. Then

$$v_{\pi}(s) = v_{\pi}^{h}(s) + v_{\pi}^{nh}(s)$$
, for $s \in S$, $\pi \in \mathcal{C}$.

Let ρ be a metric defined as in Chapter 3. Let Λ be the binary operator such that

$$x \wedge y = \min(x,y)$$
, for $x \in \mathbb{R}$, $y \in \mathbb{R}$.

We are now ready to state the following theorem.

<u>Theorem 3</u>: If π^* is an unimprovable policy such that $v_* \leq v_{\pi_0}$ and, in addition,

$$\rho(\mathbf{v}_{\pi_{O}}^{\mathrm{nh}}, \mathbf{v}_{\pi_{O}}^{\mathrm{nh}} \wedge \mathbf{v}_{\pi}^{\mathrm{nh}}) < \infty, \text{ for } \pi \in \mathcal{P},$$

then π^{\star} is optimal.

Proof: By Theorem 18 of Chapter 3, we only have to show that

$$\rho(\mathbf{v}_{\pi^*}, \mathbf{v}_{\pi^*} \wedge \mathbf{v}_{\pi}) < \infty$$
, for $\pi \in \mathcal{O}$.

But

$$\begin{split} \rho(\mathbf{v}_{\pi^{*}}, \mathbf{v}_{\pi^{*}} \wedge \mathbf{v}_{\pi}) &\leq \rho(\mathbf{v}_{\pi_{0}}, \mathbf{v}_{\pi_{0}} \wedge \mathbf{v}_{\pi}) \\ &= \rho((\mathbf{v}_{\pi_{0}}^{h} + \mathbf{v}_{\pi_{0}}^{nh}), (\mathbf{v}_{\pi_{0}}^{h} + \mathbf{v}_{\pi_{0}}^{nh}) \wedge (\mathbf{v}_{\pi}^{h} + \mathbf{v}_{\pi}^{nh})) \\ &\leq \rho((\mathbf{v}_{\pi_{0}}^{h} + \mathbf{v}_{\pi}^{nh}), (\mathbf{v}_{\pi_{0}}^{h} + \mathbf{v}_{\pi}^{nh}) \wedge (\mathbf{v}_{\pi_{0}}^{h} + \mathbf{v}_{\pi}^{nh})) \\ &\leq \rho(\mathbf{v}_{\pi_{0}}^{nh}, \mathbf{v}_{\pi_{0}}^{nh} \wedge \mathbf{v}_{\pi}^{nh}) \\ &\leq \rho(\mathbf{v}_{\pi_{0}}^{nh}, \mathbf{v}_{\pi_{0}}^{nh} \wedge \mathbf{v}_{\pi}^{nh}) \\ &\leq \infty, \quad \text{for} \quad \pi \in \mathfrak{S} . \end{split}$$

Q.E.D.

「日白湯湯「四」

Example (The M/G/l queueing system with removable server):

In this case, let $\|\cdot\|$ be the sup norm. Excluding those policies which turn the server on and off repeatedly at a decision epoch, then the contraction condition of Section 7 of Chapter 3 is satisfied. For any $\pi \in \mathbb{P}$,

$$\begin{aligned} \mathsf{p}(\mathsf{v}_{\pi_0}^{\mathrm{nh}}, \, \mathsf{v}_{\pi_0}^{\mathrm{nh}} \land \, \mathsf{v}_{\pi}^{\mathrm{nh}}) &= \|\mathsf{v}_{\pi_0}^{\mathrm{nh}} - \, \mathsf{v}_{\pi_0}^{\mathrm{nh}} \land \, \mathsf{v}_{\pi}^{\mathrm{nh}}\| \\ &= \sup_{s \in S} |\mathsf{v}_{\pi_0}^{\mathrm{nh}}(s) - \, \mathsf{v}_{\pi_0}^{\mathrm{nh}}(s) \land \, \mathsf{v}_{\pi}^{\mathrm{nh}}(s)| \\ &\leq \infty \end{aligned}$$

We conclude that if π^* is an unimprovable policy such that $v_{\pi^*} \leq v_{\pi_0}$, then π^* is optimal.

3.4 <u>Comparison with a Policy which Minimizes the Expected Discounted</u> Holding Cost until a Finite Set of States is Reached.

We now generalize the result of Section 3. This time, let π_0 denote a policy which minimizes the expected discounted holding cost incurred until a given, finite set of states is reached. Assume that $v_{\pi_0}^h$ is finite-valued. Let ρ be defined as before. <u>Theorem 4</u>: If π^* is an unimprovable policy such that $v_{\pi^*} \leq v_{\pi_0}$

and, in addition,

$$\rho(\mathbf{v}_{\pi_{O}}^{\mathrm{nh}}, \mathbf{v}_{\pi_{O}}^{\mathrm{nh}} \wedge \mathbf{v}_{\pi}^{\mathrm{nh}}) < \infty, \text{ for } \pi \in \mathcal{P},$$

then π^* is optimal.

<u>Proof</u>: Since π_0 minimizes the expected discounted holding cost incurred until a given, finite set of states is reached, and since $v_{\pi_0}^h$ is finite-valued, there must exist an $M < \infty$ such that

$$v_{\pi_0}^{h}(s) \leq v_{\pi}^{h}(s) + M$$
, for $s \in S$, $\pi \in \Theta$.

Now

$$\rho(\mathbf{v}_{\pi^*}, \mathbf{v}_{\pi^*} \wedge \mathbf{v}_{\pi}) \leq \rho(\mathbf{v}_{\pi_0}, \mathbf{v}_{\pi_0} \wedge \mathbf{v}_{\pi})$$

$$= \rho(\mathbf{v}_{\pi_0}, \mathbf{v}_{\pi_0} \wedge (\mathbf{v}_{\pi}^{\mathrm{h}} + \mathbf{v}_{\pi}^{\mathrm{nh}}))$$

$$\leq \rho(\mathbf{v}_{\pi_0}, \mathbf{v}_{\pi_0} \wedge (\mathbf{v}_{\pi_0}^{\mathrm{h}} - \mathbf{M}\theta + \mathbf{v}_{\pi}^{\mathrm{nh}}))$$

$$\leq \rho(\mathbf{v}_{\pi_0}^{\mathrm{nh}}, \mathbf{v}_{\pi_0}^{\mathrm{nh}} \wedge (\mathbf{v}_{\pi}^{\mathrm{nh}} - \mathbf{M}\theta))$$

$$\leq \rho(\mathbf{v}_{\pi_0}^{\mathrm{nh}}, \mathbf{v}_{\pi_0}^{\mathrm{nh}} \wedge \mathbf{v}_{\pi}^{\mathrm{nh}}) + \mathbf{M}$$

$$< \infty, \text{ for } \pi \in \mathfrak{S}.$$

Theorem 18 of Chapter 3 now implies that π^{\star} is optimal.

Q.E.D.

Example (The M/G/l queueing system with removable server):

Let \mathcal{J} be the set of policies such that each policy in \mathcal{J} always turns the server on (or keeps him on) when the number of customers in the system is sufficiently large. It is easy to show that, for each $\pi \in \mathcal{J}$, either $v_{\mathcal{T}}^{h}(s)$ is finite-valued for each s in S or it is infinite-valued for each s in S. In the latter case, all policies may be regarded as optimal. Therefore, we now focus on the former case where $v_{\mathcal{T}}^{h}(s)$ is always finite-valued.

Clearly, π_0 in the theorem may be any policy in $\mathcal X$. As before,

$$\rho(\mathbf{v}_{\pi_{O}}^{\mathrm{nh}}, \mathbf{v}_{\pi_{O}}^{\mathrm{nh}} \wedge \mathbf{v}_{\pi}^{\mathrm{nh}}) < \infty, \text{ for } \pi \in \mathfrak{S}.$$

Therefore, we conclude that if π^* is an unimprovable policy in $\overset{\,\,{}_{\,\!\!\!\!\!\!\!}}{$, then π^* is also optimal.

REFERENCES

Bell, C. (1971), "Characterization and Computation of Optimal Policies for Operating an M/G/1 Queueing System with Removable Server," <u>Oper. Res. 19</u>, 208-218.

Bellman, R. (1957), Dynamic Programming, Princeton University Press.

- Blackburn, J. (1971), "Optimal Control of Queueing Systems with Intermittent Service," Tech. Rep. No. 8, Department of Operations Research, Stanford University.
- Blackwell, D. (1962), "Discrete Dynamic Programming," <u>Ann. Math. Stat.</u> <u>33</u>, <u>7</u>19-726.
- Blackwell, D. (1965), "Discounted Dynamic Programming," <u>Ann. Math. Stat.</u> <u>36</u>, 226-235.

Dantzig, G. B. and Wolfe, P. (1962), "Linear Programming in a Markov Chain," Oper. Res. 10, 707-710.

- Denardo, E. (1967), "Contraction Mappings in the Theory Underlying Dynamic Programming," <u>SIAM Rev. 9</u>, 165-177.
- Denardo, E. V. (1970a), "On Linear Programming in a Markov Decision Problem," <u>Mgt. Sci.</u> 16, 281-288.

Denardo, E. V. (1970b), "Computing Bias-Optimal Policies in Discrete and Continuous Markov Decision Problems," Oper. <u>Res. 18</u>, 279-289.

- Denardo, E. V. (1971), "Markov Renewal Programs with Small Interest Rates," Ann. Math. Stat. 42, No. 2, 477-496.
- Denardo, E. V. and Fox, B. L. (1968), "Multichain Markov Renewal Programs," SIAM J. Appl. Math. 16, 468-487.

- Derman, C. (1966), "Denumerable State Markovian Decision Processes -Average Cost Criterion," Ann. Math. Stat. 37, 1545-1554.
- Derman, C. (1970), <u>Finite State Markovian Decision Processes</u>, Academic Press.
- D'Epenoux, F. (1960), "Sur un Probleme de Production et de Stockage dans l'a Leatoire," <u>Rev. Francaise Informat. Recherche Oper-</u> <u>ationelle 14</u>, 3-16 [English Transl.: <u>Mgt. Sci. 10</u>, 98-108 (1963).]
- Gross, D. and Harris, C. M. (1974), <u>Fundamentals of Queueing Theory</u>, John Wiley and Sons, Inc.
- Harrison, M. (1972), "Discrete Dynamic Programming with Unbounded Rewards," Ann. Math. Stat. 43, 636-644.
- Heyman, D. (1968), "Optimal Operating Policies for M/G/l Queueing Systems," Oper. Res. 16, 362-382.
- Hordijk, A. (1974a), <u>Dynamic Programming and Markov Potential Theory</u>, Matematisch Centrum.
- Hordijk, A. (1974b), Jonvergent Dynamic Programming," Tech. Rep. No. 28, Department of Operations Research, Stanford University.
- Howard, R. (1960), <u>Dynamic Programming and Markov Processes</u>, Technology Press of M.I.T., Cambridge.
- Howard, R. A. (1964), "Research in Semi-Markovian Decision Structure," J. Oper. Res. Soc. Japan 6, No. 4.
- Iglehart, D. L. (1963), "Optimality of (s,S) Policies in the Infinite Horizon Dynamic Inventory Problem," Mgt. Sci. 9, 259-267.
- Jewell, W. S. (1963), "Markov Renewal Programming, I and II," Oper. Res. 11, 938-971.

Lippman, S. (1975a), "On Dynamic Programming with Unbounded Rewards,"

Mgt. Sci. 27, 1225-1233.

Lippman, S. A. (1976a), "Applying a New Device in the Optimization of Exponential Queueing Systems, <u>Oper. Res. 23</u>, 687-710.

Magazine, M. (1971), "Optimal Control of Multi-Channel Service Systems," Nav. Res. Log. Quart. 18, 177-183. Manne, A. (1960), "Linear Programming and Sequential Decisions," Mgt. Sci. 6, No. 3, 259-267.

- McGill, J. T. (1969), "Optimal Control of Queueing Systems with Variable Number of Exponential Servers," Tech. Rep. No. 123, Department of Operations Research, Stanford University.
- Orkenyi, P. (1976), "Optimal Control of the M/G/l Queueing System with Removable Server - Linear and Non-Linear Holding Cost Function," Tech. Rep. No. 65, Department of Operations Research, Stanford University. Office of Naval Research Contract N00014-76-C-0418.

Porteus, E. L. (1975a), "An Informal Look at Principle of Optimality," <u>Mgt. Sci.</u> 21, 1346-1348.

- Porteus, E. L. (1974b), "On the Optimality of Structured Policies in Countable Stage Decision Processes," <u>Mgt. Sci</u>. <u>22</u>, 148-158.
- Prabhu, N. U. and Stidham, Jr. S. (1973), "Optimal Control of Queueing Systems," in Mathematical Methods in Queueing Theory, Conference at Western Michigan University, May 10-12.
- Pyke, R., "Markov Renewal Processes with Finitely Many States," <u>Ann</u>. Math. Stat. <u>3</u>2, 1243-1259.
- Reed, C. (1973), "Denumerable State Decision Processes with Unbounded Costs," Tech. Rep. No. 22, Department of Operations Research, Stanford University.

77

- Reed, C. (1974a), "Difference Equations and the Optimal Control of Single Server Queueing Systems," Tech. Rep. No. 23, Department of Operations Research, Stanford University.
- Reed, F. C. (1974b), "The Effect of Stochastic Time Delays on Optimal Operating Policies for M/G/l Queueing Systems with Intermittent Service," Tech. Rep. No. 45, Department of Operations Research, Stanford University.
- Ross, S. (1968), "Arbitrary State Markovian Decision Processes," <u>Ann</u>. <u>Math. Stat. 39</u>, 2118-2122.
- Ross, S. (1970), <u>Applied Probability Models with Optimization Appli-</u> cations, Holden-Day.
- Smith, W. L. (1955), "Regenerative Stochastic Processes," <u>Proceedings</u> <u>Royal Society</u>, Series A, 232, 6-31.
- Strauch, R. E. (1966), "Negative Dynamic Programming," <u>Ann. Math. Stat.</u> <u>37</u>, 871-890.
- Veinott, A. F. Jr. (1966), "On Finding Optimal Policies in Discrete
 Dynamic Programming with No Discounting," <u>Ann. Math. Stat. 37</u>,
 1284-1294.
- Veinott, A. F. Jr. (1969), "Discrete Dynamic Programming with Sensitive Discount Optimality Criteria," <u>Ann. Math. Stat.</u> <u>40</u>, 1635-1660.

78

ECURITY CLASSIFICATION OF THIS PAGE (When Date 2	Intered)	
REPORT DOCUMENTATION	PAGE	READ INSTRUCTIONS BEFORE COMPLETING FORM
TR 464	2. GOVT ACCESSION NO.	3. BESIPIENT'S CATALOG NUMBER
A TITLE (and Subilite)		S. THE OF REPORT & PERIOD COVERED
UNBOUNDED COSTS AND ITS APPLICATION	ON TO THE	Technical Repert,
OPTIMAL CONTROL OF QUEUEING SYSTEM	s,	6. FERFORMING ORG. REPORT NUMBER
7. AUTHOR(a)	<u> </u>	. CONTRACT OR GRANT NURBER(.)
() Peter Orkenyi		N00014-76-C-0418
2		VNSF-5ng-75-118
Dept. of Operations Research	× 277	10. PROGRAM ELEMENT, PROJECT, TASK
Stanford University Stanford, California	$2 \delta_{\alpha} \rho_{i}$	26-(NR-\$47-061)
1. CONTROLLING OFFICE NAME AND ADDRESS	ce Branch	A REPORT DATE
Code 434	co branch	11 August 976
Office of Naval Research - Arling	ton, Va. 22217	78
A A A A A A A A A A A A A A A A A A A		Inclassified
		The DECI AMERICANION PERMETANI
		SCHEDULE
		ited.
7. DISTRIBUTION STATEMENT (of the abelract entered in	Block 20, If different free	ited. m Report)
7. DISTRIBUTION STATEMENT (of the ebstract entered in 8. SUPPLEMENTARY NOTES	Block 20, If different fro	ited. m Repotij
 DISTRIBUTION STATEMENT (of the ebstrect entered in s. SUPPLEMENTARY NOTES This research was supported in p Grant ENG 75-14847 and The Norw the Humanities. 	Block 20, If different fro part by National egian Research (ited. Report) I Science Foundation Council for Science and
 DISTRIBUTION STATEMENT (of the ebstract entered in S. SUPPLEMENTARY NOTES This research was supported in Grant ENG 75-14847 and The Norw the Humanities. KEY WORDS (Continue on reverse side if necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV 	part by Nationa Identify by block number) DECISION PROCE	ited. Report) I Science Foundation Council for Science and SSES, QUEUEING THEORY,
 DISTRIBUTION STATEMENT (of the ebstrect entered in particular statement of the ebstrect entered in particular statement was supported in part ENG 75-14847 and The Norwethe Humanities. KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE 1 	part by Nationa egian Research (Identify by block number) DECISION PROCE POLICIES, OPTIN	Ited. Report Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS,
 7. DISTRIBUTION STATEMENT (of the ebstrect entered in particular statement of the ebstrect entered in particular statement was supported in part ENG 75-14847 and The Norwathe Humanities. A KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE POLICY IMPROVEMENT, STATIONARY OF 	part by National egian Research (Identify by block number) DECISION PROCES POLICIES, OPTII PTIMLL POLICIES	Ited. Report Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS,
 DISTRIBUTION STATEMENT (of the ebstrect entered in B. SUPPLEMENTARY NOTES This research was supported in) Grant ENG 75-14847 and The Norwe the Humanities. KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE T POLICY IMPROVEMENT, STATIONARY OF ABSTRACT (Continue on reverse elde II necessary and I 	part by National egian Research (Identify by block number) DECISION PROCE POLICIES, OPTIM PTIMAL POLICIES	Ited. Report L Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS,
 DISTRIBUTION STATEMENT (of the ebstrect entered in SUPPLEMENTARY NOTES This research was supported in : Grant ENG 75-14847 and The Norwe the Humanities. KEY WORDS (Continue on reverse eldo if necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE : POLICY IMPROVEMENT, STATIONARY O: ABSTRACT (Continue on reverse elde if necessary and in 	part by Nationa egian Research (Identify by block number) DECISION PROCES POLICIES, OPTIN PTIMIL POLICIES	Ited. Report SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered in B. SUPPLEMENTARY NOTES This research was supported in g Grant ENG 75-14847 and The Norw the Humanities. 7. KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE T POLICY IMPROVEMENT, STATIONARY OF ABSTRACT (Continue on reverse elde II necessary and in 	Block 20, If different from part by National egian Research (Identify by block number) DECISION PROCE POLICIES, OPTIM PTIMAL POLICIES	TReport)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in grant ENG 75-14847 and The Norw, the Humanities. 7. KEY WORDS (Continue on reverse elde if necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE POLICY IMPROVEMENT, STATIONARY OPPLICY IMPROVEMENT, STATIONARY IMPROVEMENT, STATIONARY	part by National egian Research (Identify by block number) DECISION PROCES POLICIES, OPTIM PTIMAL POLICIES	Ited. Report) L Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in : Grant ENG 75-14847 and The Norw the Humanities. 9. KEY WORDS (Continue on reverse eldo if necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE : POLICY IMPROVEMENT, STATIONARY O: ABSTRACT (Continue on reverse elde if necessary and is	part by Nationa egian Research (Identify by block number) DECISION PROCES POLICIES, OPTIN PTIMIL POLICIES	ited. Report Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in grant ENG 75-14847 and The Norw, the Humanities. 7. KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE POLICY IMPROVEMENT, STATIONARY OF SUPPLICY IMPROVEMENT, STATIONARY IMPROVEMENT, STATIONARY OF SUPPLICY IMPROVEMENT, SUPPLICY IMPROVEMENT,	part by National egian Research (Identify by block number) DECISION PROCES POLICIES, OPTIM PTIMAL POLICIES	Report) A Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in grant ENG 75-14847 and The Norw the Humanities. 9. KEY WORDS (Continue on reverse eldo II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE POLICY IMPROVEMENT, STATIONARY O. 9. ABSTRACT (Continue on reverse elde II necessary and I taken 1475) 	part by National egian Research (Identify by block number) DECISION PROCES POLICIES, OPTIN PTIMAL POLICIES	ited. Report: Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side)
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in ; Grant ENG 75-14847 and The Norw, the Humanities. 9. KEY WORDS (Continue on reverse eldo if necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE : POLICY IMPROVEMENT, STATIONARY O 9. ABSTRACT (Continue on reverse elde if necessary and i ABSTRACT (Continue on reverse elde if necessary and i O 1 JAN 73 1475 EDITION OF 1 NOV 65 15 OBSOLE S/M 0102-014-6601 1 	Block 20, if different from part by National egian Research (Identify by block number) DECISION PROCE: POLICIES, OPTIN PTIMIL POLICIES (dentify by block number)	Ited. Report Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side) NCLASSIFIED
 7. DISTRIBUTION STATEMENT (of the ebstrect entered it 8. SUPPLEMENTARY NOTES This research was supported in grant ENG 75-14847 and The Norw the Humanities. 9. KEY WORDS (Continue on reverse elde II necessary and DYNAMIC PROGRAMMING, SEMI-MARKOV QUEUEING SYSTEMS, UNIMPROVABLE POLICY IMPROVEMENT, STATIONARY OPPLICY I JAN 73 1475 EDITION OPPLICY IMPROVEMENT, SOUTH STATIONARY SISONALE S/M 0102-014- 6601 I	TE	Ited. Report) I Science Foundation Council for Science and SSES, QUEUEING THEORY, MALITY CONDITIONS, (see reverse side) NCLASSIFIED SEFICATION OF THIS PAGE (From Data Entere

Ġ

:

UNCLASSIFIED SECURITY CLASSIFICATION OF THIS PAGE (When Date Enjoyed)

> A THEORY FOR SEMI-MARKOV DECISION PROCESSES WITH UNBOUNDED COSTS AND ITS APPLICATION TO THE OPTIMAL CONTROL OF QUEUEING SYSTEMS

> > Ъy

Peter Orkenyi

Abstract:

2 Semi-Markov decision processes with countable state and action spaces are investigated. The optimality criteria considered are the average cost criterion, the undiscounted cost criterion, and the discounted cost criterion. The common assumption of bounded costs has been replaced by some considerably weaker conditions. In particular, our assumptions are weaker than those made by Harrison, Hordijk, Lippman and Reed when they considered the same problem.

The existence of optimal, stationary optimal and stationary \mathcal{E} -optimal policies is investigated. Policy improvement is considered. Necessary and sufficient conditions for the optimality of a policy are given.

Then the optimal control of queueing systems is considered by formulating this general problem as a semi-Markov decision process. Finally, four different ways of proving the optimality of an unimprovable policy are developed in the context of queueing systems.

UNCLASSIFIED SECURITY CLASSIFICATION OF THIS PAGE(When Deta Entered)