

DISCRETE-TIME CONTROLLED MARKOV PROCESSES WITH AVERAGE COST CRITERION: A SURVEY*

ARISTOTLE ARAPOSTATHIS[†], VIVEK S. BORKAR[‡], EMMANUEL
FERNÁNDEZ-GAUCHERAND[§], MRINAL K. GHOSH¹, AND STEVEN I. MARCUS²

This paper is dedicated to Wendell Fleming on the occasion of his 65th birthday.

Abstract. This work is a survey of the average cost control problem for discrete-time Markov processes. The authors have attempted to put together a comprehensive account of the considerable research on this problem over the past three decades. The exposition ranges from finite to Borel state and action spaces and includes a variety of methodologies to find and characterize optimal policies. The authors have included a brief historical perspective of the research efforts in this area and have compiled a substantial yet not exhaustive bibliography. The authors have also identified several important questions that are still open to investigation.

Key words. controlled Markov processes, average cost, stationary policies, dynamic programming, optimal policies, ergodicity

AMS(MOS) subject classifications. 93E20, 60J70

1. Introduction. The average cost criterion (equivalently, the long-run average or ergodic cost) is a popular criterion for optimization of stochastic dynamical systems over an infinite time horizon. It is a reasonable criterion to use when the anticipated time interval for optimization (which in practice is finite) is long compared to other timescales involved, and there are no compelling reasons to prefer short-term optimization over long-term. Naturally, it is not favored in financial applications where money spent now is worth more than money spent later, but there are situations (communication networks being a prime example) where a “steady state” operation is expected over intervals that are long compared to the time constants of the system. Then it makes sense to minimize the limiting time-averaged cost, i.e., the “average cost.”

Mathematically, the criterion stands out as being much more difficult to analyze than the others; while other classical criteria lead to reasonably complete solutions, the average cost does not. The finite state and action problem is well understood, but there are numerous counterexamples in which infinite state or action problems do not have a nice solution. In fact, it appears not as a single problem but a collection of problems, some of which do not have a nice solution (cf. [150]). Thus, a variety of approaches have been developed to handle different situations. Not surprisingly,

*Received by the editors October 25, 1991; accepted for publication (in revised form) July 21, 1992. This work was supported in part by Texas Advanced Research Program (Advanced Technology Program) grants 003658-093 and 003658-186; in part by Air Force Office of Scientific Research grants AFOSR-91-0033, F49620-92-J-0045, and F49620-92-J-0083; and in part by National Science Foundation grant CDR-8803012.

[†]Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin, Texas 78712.

[‡]Department of Electrical Engineering, Indian Institute of Science, Bangalore 560012, India.

[§]Systems and Industrial Engineering Department, University of Arizona, Tucson, Arizona 85721.

¹Department of Mathematics, Indian Institute of Science, Bangalore 560012, India.

²Department of Electrical Engineering and Systems Research Center, University of Maryland, College Park, Maryland 20742.

this is one chapter of Markov decision theory that is anything but closed. At the same time, it has come of age, having been studied for over 30 years, with promises of significant advances on the horizon. This, in short, is the *raison d'être* for this survey; we have attempted to put together a coherent account of what has been done, with an indication of what future advances may be.

Any such project has obvious limitations. Space constraints dictate a certain amount of selection, and not every relevant work can be covered in significant detail. We have included proofs where we felt they were essential to understanding the results or contained potentially useful novel ideas. In all cases, a serious attempt at objectivity has been made. For complementary reading on the general subject of Markov decision theory, see [137], [181], [196], [207].

The paper is organized as follows: §2 describes the problem formulation in full detail. Section 3 gives a brief history. Sections 4–6 extensively treat the finite state, the countable state, and the Borel state space cases, respectively, under complete observations. Section 7 treats the problem under partial observations. Section 8 describes some recent results on multiobjective average cost control. Finally, we conclude with some relevant remarks.

2. Preliminaries and formulation of the problem. In this section, the model and basic results concerning controlled Markov processes are given in the most general form needed for our presentation. In some subsequent sections, we specialize our presentation to situations in which measure-theoretic aspects are of no essential concern, as in the case for models with countable state space, allowing for a more transparent exposition. Before presenting the model, we summarize our key notation as follows:

- \mathbb{R} : set of real numbers;
- \mathbb{N} : set of positive integers;
- \mathbb{N}_0 : set of nonnegative integers;
- $\mathcal{B}(\mathbf{W})$: Borel σ -algebra of a given topological space \mathbf{W} ;
- $\mathcal{P}(\mathbf{W})$; for a Borel space \mathbf{W} (see [15], [82]), the set of all probability measures on $\mathcal{B}(\mathbf{W})$ endowed with the topology of weak convergence (see [134]).

The following are function spaces on a topological space \mathbf{W} :

- $C_b(\mathbf{W}) := \{v : \mathbf{W} \rightarrow \mathbb{R} \mid v \text{ is continuous and bounded}\}$;
- $\mathcal{M}(\mathbf{W}) := \{v : \mathbf{W} \rightarrow \mathbb{R} \mid v \text{ is Borel measurable}\}$;
- $\mathcal{M}_b(\mathbf{W}) := \{v : \mathbf{W} \rightarrow \mathbb{R} \mid v \text{ is Borel measurable and bounded}\}$;
- $\mathcal{L}(\mathbf{W}) := \{v : \mathbf{W} \rightarrow \mathbb{R} \mid v \text{ is lower semicontinuous and bounded below}\}$;
- $\mathcal{L}_b(\mathbf{W}) := \mathcal{L}(\mathbf{W}) \cap \mathcal{M}_b(\mathbf{W})$.

For $v \in \mathcal{M}_b(\mathbf{W})$, we let

- $\|v\| := \sup_{w \in \mathbf{W}} \{|v(w)|\}$;
- $\text{span}(v) := \sup_{w, w' \in \mathbf{W}} \{v(w) - v(w')\}$;
- $v^+ := v - \inf_{w \in \mathbf{W}} \{v(w)\}$, $v^- := v - \sup_{w \in \mathbf{W}} \{v(w)\}$.

We refer to $\text{span}(v)$ as the *span seminorm* of v .

The following is a list of the abbreviations used in this paper (the section where each abbreviation is first introduced is indicated in parenthesis):

- AC average cost (§2.4);
- ACOE average cost optimality equation (§3);
- ACOI average cost optimality inequality (§5.2);
- CMP controlled Markov process (§2.1);

- CO completely observable (§3);
- DC discounted cost (§2.4);
- DCOE discounted cost optimality equation (§2.6);
- PO partially observable (§7.2);
- POCMP partially observable controlled Markov process (§3);
- TC total cost (§2.4).

2.1. The model. A discrete-time, stationary controlled Markov process (CMP), or Markov decision process, is a stochastic dynamical system specified by the five-tuple $(\mathbf{S}, \mathbf{A}, U, P, c)$, where

(a) \mathbf{S} is a Borel space, called the *state space*, the elements of which are called *states*;

(b) \mathbf{A} is a Borel space, called the *action* or *control* space;

(c) $U : \mathbf{S} \rightarrow \mathcal{B}(\mathbf{A})$ is a strict, measurable, compact-valued multifunction (see the Appendix). $U(x)$ represents the set of admissible actions (or control inputs) when the system is in state $x \in \mathbf{S}$. Accordingly, the set of admissible state/action pairs is $\mathbf{K} := \{(x, a) : x \in \mathbf{S}, a \in U(x)\} = \text{Graph}(U)$, and we have that $\mathbf{K} \in \mathcal{B}(\mathbf{S} \times \mathbf{A})$. This set is endowed with the subspace topology corresponding to $\mathcal{B}(\mathbf{S} \times \mathbf{A})$;

(d) P is a stochastic kernel on \mathbf{S} given \mathbf{K} , called the *transition kernel*. It is assumed to be Borel measurable, i.e., $P(D | \cdot) : \mathbf{K} \rightarrow [0, 1]$ is Borel measurable, for each $D \in \mathcal{B}(\mathbf{S})$;

(e) $c : \mathbf{K} \rightarrow \mathbb{R}$ is the (measurable) one-stage cost function.

The evolution of the system is as follows. Let X_t denote the state at time $t \in \mathbb{N}_0$, and A_t the action chosen at that time. If $X_t = x \in \mathbf{S}$ and $A_t = a \in U(x)$, then (i) a cost $c(x, a)$ is incurred, and (ii) the system moves to the next state X_{t+1} , according to a probability distribution $P(\cdot | x, a)$. Once the transition into the next state has occurred, a new action is chosen, and the process is repeated.

The total period of time over which the system is to be observed is called the planning (or decision-making or control) horizon and is denoted by T . It can be a finite interval $\{0, \dots, N - 1\}$, with $N \in \mathbb{N}$, or an infinite horizon, e.g., $T = \mathbb{N}_0$.

The (admissible) *history spaces* are defined as

$$\mathbf{H}_0 := \mathbf{S}, \quad \mathbf{H}_t := \mathbf{H}_{t-1} \times \mathbf{K}, \quad t \in \mathbb{N}_0,$$

and the canonical sample space is defined as $\mathbf{\Omega} := (\mathbf{S} \times \mathbf{A})^\infty$. These spaces are endowed with their respective product topologies and are therefore Borel spaces. A generic element $\omega \in \mathbf{\Omega}$ is of the form $\omega = (x_0, a_0, x_1, a_1, \dots)$, $x_i \in \mathbf{S}$, $a_i \in \mathbf{A}$; all random variables will be defined on the measurable space $(\mathbf{\Omega}, \mathcal{B}(\mathbf{\Omega}))$.

The state, action (or control), and information processes, denoted by $\{X_t\}_{t \in T}$, $\{A_t\}_{t \in T}$ and $\{H_t\}_{t \in T}$, respectively, are defined by the projections

$$X_t(\omega) := x_t, \quad A_t(\omega) := a_t, \quad H_t(\omega) := (x_0, \dots, a_{t-1}, x_t), \quad t \in T$$

for each realization $\omega = (x_0, \dots, a_{t-1}, x_t, a_t, \dots) \in \mathbf{\Omega}$. Since $\mathcal{B}(\mathbf{\Omega}) = (\mathcal{B}(\mathbf{S}) \times \mathcal{B}(\mathbf{A}))^\infty$, the above are well-defined random processes on $(\mathbf{\Omega}, \mathcal{B}(\mathbf{\Omega}))$. Note that $\mathcal{B}(\mathbf{\Omega}) = \bigvee_{t=0}^\infty \mathfrak{F}_t$, where $\mathfrak{F}_t = \sigma(H_t)$, the σ -algebra generated by H_t .

Example 2.1. Let \mathbf{S} , \mathbf{A} , \mathbf{W} be Borel spaces and $F : \mathbf{S} \times \mathbf{A} \times \mathbf{W} \rightarrow \mathbf{S}$ a Borel function. Consider a nonlinear stochastic system described by the system equation

$$X_{t+1} = F(X_t, A_t, W_t), \quad t \in T,$$

where the process $\{W_t\}$ is a sequence of independent and identically distributed (i.i.d.) \mathbf{W} -valued random variables, with common probability distribution \mathcal{P}_W , often referred to as a stochastic state disturbance, or noise; $\{W_t\}$ is assumed to be independent of X_0 . Suppose that a strict, measurable, compact-valued multifunction $U : \mathbf{S} \rightarrow \mathcal{B}(\mathbf{A})$ has been specified, giving the necessary constraints on the control actions, or that $U(x) = \mathbf{A}$, for all $x \in \mathbf{S}$, if there are no constraints. Then the evolution of the system is equivalently described in terms of the stochastic kernel P on \mathbf{S} given \mathbf{K} defined as

$$P(D \mid x, a) := \int_{\mathbf{W}} I\{F(x, a, w) \in D\} \mathcal{P}_W(dw), \quad (x, a) \in \mathbf{K}, \quad D \in \mathcal{B}(\mathbf{S}),$$

where $I\{A\}$ denotes the indicator function of the event A . The additional specification of a measurable cost function $c : \mathbf{K} \rightarrow \mathbb{R}$ would completely define a CMP $(\mathbf{S}, \mathbf{A}, U, P, c)$.

Example 2.2. Consider a countable set \mathbf{S} endowed with the discrete topology. With no loss in generality we can take $\mathbf{S} = \mathbb{N}_0$. Let \mathbf{A} be a Borel space and $U(x) = \mathbf{A}$, for all $x \in \mathbf{S}$. In this case, every stochastic kernel on \mathbb{N}_0 given $\mathbf{K} := \mathbb{N}_0 \times \mathbf{A}$ reduces to a collection of discrete probability distributions parameterized by $(i, a) \in \mathbf{K}$. These can also be represented by a collection of stochastic matrices $\{P(a) = [p_{ij}(a)]\}_{a \in \mathbf{A}}$; i.e., $P(a)$ is a state transition matrix, and $p_{ij}(a)$ is the probability that the state of the system makes a transition from i to j , under action a . Therefore, additionally specifying a cost function $c : \mathbb{N}_0 \times \mathbf{A} \rightarrow \mathbb{R}$ completely defines a CMP.

2.2. Policies and performance criteria. An *admissible control strategy*, or *policy*, is a sequence $\pi = \{\pi_t\}_{t \in T}$ of Borel measurable stochastic kernels on \mathbf{A} given \mathbf{H}_t , satisfying the constraint

$$\pi_t(U(x_t) \mid h_t) = 1, \quad x_t \in \mathbf{S}, \quad h_t \in \mathbf{H}_t.$$

The set of all admissible policies will be denoted by Π .

If $\mu \in \mathcal{P}(\mathbf{S})$ and $\pi \in \Pi$ are given, there exists a unique probability measure \mathcal{P}_μ^π on $(\Omega, \mathcal{B}(\Omega))$ satisfying the following [15, Prop. 7.28, pp. 140–144], [130, Prop. V.1.1, pp. 162–164], with $D \in \mathcal{B}(\mathbf{S})$ and $C \in \mathcal{B}(\mathbf{A})$:

$$(2.1) \quad \mathcal{P}_\mu^\pi(X_0 \in D) = \mu(D),$$

$$(2.2) \quad \mathcal{P}_\mu^\pi(A_t \in C \mid H_t) = \pi_t(C \mid H_t), \quad \mathcal{P}_\mu^\pi\text{-a.s.},$$

$$(2.3) \quad \mathcal{P}_\mu^\pi(X_{t+1} \in D \mid H_t, A_t) = P(D \mid X_t, A_t), \quad \mathcal{P}_\mu^\pi\text{-a.s.}$$

Therefore, if μ is the distribution of the initial state X_0 , and policy $\pi \in \Pi$ is used, the underlying probability space of all random variables of interest is $(\Omega, \mathcal{B}(\Omega), \mathcal{P}_\mu^\pi)$. The expectation operator with respect to \mathcal{P}_μ^π will be denoted by E_μ^π . Furthermore, if μ is a Dirac measure at $x \in \mathbf{S}$, we will simply write \mathcal{P}_x^π and E_x^π .

Certain classes of admissible policies are of special interest. A policy π is called a *Markov randomized policy* if there exists a sequence of measurable maps $\{f_t\}_{t \in T}$, called *randomized decision rules*, where $f_t : \mathbf{S} \rightarrow \mathcal{P}(\mathbf{A})$, for each $t \in T$, such that

$$\pi_t(\cdot \mid H_t) = f_t(X_t)(\cdot), \quad \mathcal{P}_\mu^\pi\text{-a.s.}$$

Conversely, every sequence of measurable maps $f_t : \mathbf{S} \rightarrow \mathcal{P}(\mathbf{A})$, $t \in T$, satisfying $f_t(x)(U(x)) = 1$, defines a Markov randomized policy in an obvious way; with some

abuse in notation, the sequence itself will be referred to as the policy. The set of all Markov randomized policies will be denoted by Π_M . A policy $\{f_t\}_{t \in T} \in \Pi_M$ is called a *stationary* randomized policy if there is a randomized decision rule f such that, for all $t \in T$, $f_t = f$. The set of all stationary randomized policies will be denoted by Π_{SR} . A *nonrandomized, deterministic, or pure* decision rule is a measurable map $f : \mathbf{S} \rightarrow \mathbf{A}$. A policy $\{f_t\}_{t \in T} \in \Pi_M$ is called a *nonrandomized, deterministic, or pure* Markov policy if each f_t is deterministic. Hence, in this case, $A_t = f_t(X_t)$ almost surely. The set of deterministic Markov policies will be denoted by Π_{MD} . Stationary deterministic policies are defined in the obvious way. The set of all stationary deterministic policies is denoted by Π_{SD} , and, for $\pi \in \Pi_{SD}$, $\pi(x)$ will denote the action chosen at $x \in \mathbf{S}$. Clearly $\Pi_{SD} \subseteq \Pi_{MD} \subseteq \Pi_M \subseteq \Pi$, and $\Pi_{SD} \subseteq \Pi_{SR} \subseteq \Pi_M$.

It is easily seen that, under a policy $\pi = \{f_t\}_{t \in T} \in \Pi_M$, the state process $\{X_t\}_{t \in T}$ is a Markov process. That is, for $D \in \mathcal{B}(\mathbf{S})$,

$$\begin{aligned} \mathcal{P}_\mu^\pi(X_{t+1} \in D \mid X_t, \dots, X_0) &= \mathcal{P}_\mu^\pi(X_{t+1} \in D \mid X_t) \\ &= \int_{\mathbf{A}} P(D \mid X_t, a) f_t(X_t)(da), \quad \mathcal{P}_\mu^\pi\text{-a.s.}, \end{aligned}$$

and, under a policy $\pi' \in \Pi_{SR}$, $\{X_t\}_{t \in T}$ is a Markov process with stationary transition probabilities.

Each policy $\pi \in \Pi$ incurs a stream of random costs, e.g., $\{c(X_t, f_t(X_t))\}_{t \in T}$, for $\{f_t\}_{t \in T} \in \Pi_{MD}$. Depending upon the problem requirements, several cost evaluation criteria are studied. The following criteria are frequently used.

Total cost (TC). The total cost incurred by the policy $\pi \in \Pi$ over the entire planning horizon is given by

$$J_T(\mu, \pi) := E_\mu^\pi \left[\sum_{t \in T} c(X_t, A_t) \right].$$

When the horizon is finite, i.e., $T = \{0, \dots, N - 1\}$, $N \in \mathbb{N}_0$, we denote the above more explicitly as $J_N(\mu, \pi)$. Furthermore, given a *terminal cost* function $h \in \mathcal{M}_b(\mathbf{S})$, we define

$$J_N(\mu, \pi, h) := E_\mu^\pi \left[\sum_{t=0}^{N-1} c(X_t, A_t) + h(X_N) \right].$$

Discounted cost (DC). Let $0 < \beta < 1$, the *discount factor*, and $\pi \in \Pi$ be given. The total discounted cost incurred by π over the infinite planning horizon is given by

$$J_\beta(\mu, \pi) := E_\mu^\pi \left[\sum_{t=0}^{\infty} \beta^t c(X_t, A_t) \right].$$

Average cost (AC). The expected long-run average cost incurred by $\pi \in \Pi$ is given by

$$J(\mu, \pi) := \limsup_{N \rightarrow \infty} E_\mu^\pi \left[\frac{1}{N} \sum_{t=0}^{N-1} c(X_t, A_t) \right] = \limsup_{N \rightarrow \infty} \frac{1}{N} J_N(\mu, \pi).$$

Sample path average cost. This is a pathwise version of the AC, and, for $X_0 = x$, it is given by

$$J_S(x, \pi) := \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} c(X_t, A_t),$$

where $\{X_t\}$ and $\{A_t\}$ are the state and control process induced by $\pi \in \Pi$. Here, $J_S(x, \pi)$ is to be regarded as an extended real-valued random variable on the canonical sample space.

For the AC criterion, the limit of the expected average cost may not exist for some or all policies $\pi \in \Pi$, and thus the limit superior is used. This is always well defined and captures the *worst* possible asymptotic expected average performance under policy $\pi \in \Pi$; i.e., it gives a “pessimistic” measure of performance. On the other hand, the limit inferior could also be used, which would yield an “optimistic” measure of performance by capturing the *best* possible asymptotic expected average performance. The planning horizon for the TC criterion can be finite or infinite, whereas, for the other criteria above, it is always infinite. Under certain conditions, it can be shown that a problem with the DC criterion is equivalent to one with a TC criterion, with a random (finite) horizon; see [40, pp. 31–32]. Also, it can be shown that, for each $\pi \in \Pi$, a policy $\pi' \in \Pi_M$ can be found such that $E_\mu^\pi [c(X_t, A_t)] = E_\mu^{\pi'} [c(X_t, A_t)]$, for each $t \in \mathbb{N}_0$ and any initial distribution $\mu \in \mathcal{P}(\mathcal{S})$ [42], [51, §3.8]. Thus, for criteria that are determined by these expected costs, such as the AC, DC, and TC criteria, it suffices to consider policies in Π_M .

For an infinite planning horizon, $J_T(\mu, \pi)$ need not be well defined or may be infinite for all $\pi \in \Pi$, rendering this criterion useless for comparing the performance under different policies. Therefore, the DC or AC criteria are usually selected when the planning horizon is infinite. When the DC criterion is used, a rather complete theory is available for the corresponding dynamic programming formulation of the problem [14], [15], [51], [82], [103], [150], [200]. In this situation, future costs are discounted at a fixed rate $0 < \beta < 1$, and therefore, if β is not sufficiently close to 1, the asymptotic behavior of the state/cost process may not be important at all. Quite the opposite is the case with the AC criterion, under which all decision times are given equal weight, and we take the limit of time-averaged expected costs. The finite time evolution of the state/cost process is, in some sense, completely irrelevant in this case, and some sort of asymptotic stable behavior is desired, making this case mathematically much more involved than the previous one. Hence, the DC and AC can be seen as two opposite extremes in the spectrum of possible criteria that can be considered, in the sense that the first one captures primarily the performance of the process at the present and near future, and the second captures the performance at the distant future.

2.3. The optimal control problem. The *optimal control (or decision) problem* is that of selecting an admissible policy such that a given performance criterion is minimized over all admissible policies. For example, for the DC criterion, a policy $\pi^* \in \Pi$ is said to be (β) -discount ε -optimal for the initial distribution μ if

$$J_\beta(\mu, \pi^*) \leq J_\beta(\mu, \pi) + \varepsilon \quad \forall \pi \in \Pi,$$

where $\varepsilon > 0$. If a policy is discount ε -optimal for all distributions $\mu \in \mathcal{P}(\mathcal{S})$, then it is simply called discount ε -optimal. If a policy is discount ε -optimal *for all* $\varepsilon > 0$, then it is called *discount optimal*. The (optimal) value function is given by

$$(2.4) \quad J_\beta^*(\mu) := \inf_{\pi \in \Pi} J_\beta(\mu, \pi).$$

Also, if μ is concentrated at $x \in \mathcal{S}$, we denote the value function by $J_\beta^*(x)$. Similar definitions apply to other criteria; $J_T^*(\mu)$ and $J^*(\mu)$ will denote the optimal value functions for the TC and AC criteria, respectively. For sample path AC, we define an

optimal policy as follows: We say that a policy $\pi^* \in \Pi$ is *sample path AC optimal* (or *almost surely AC optimal*) if there exists a constant ρ^* such that, for any initial law μ ,

$$J_S^*(\mu, \pi^*) = \rho^*, \quad \mathcal{P}_\mu^{\pi^*}\text{-a.s.},$$

while, for any other policy $\pi \in \Pi$ and any initial law μ' ,

$$J_S^*(\mu', \pi) \geq \rho^*, \quad \mathcal{P}_{\mu'}^\pi\text{-a.s.}$$

The constant ρ^* is the sample path optimal average cost.

Having defined various optimality criteria and the set of admissible policies Π , the obvious question now is: Do there exist optimal policies? Without imposing further assumptions on our general model, the answer is no. One of the reasons behind this is that the Borel measurability assumption in the definition of admissible policies is too restrictive, in general, to be able to attain the infimum in (2.4). To circumvent this problem, either a broader sense of measurability is allowed, i.e., a larger set of admissible policies is used, or further assumptions are imposed. The first approach was taken by Shreve and Bertsekas [15], [164], [165], who considered *universally measurable* policies, a class properly containing the (Borel measurable) admissible policies defined previously; see also [51]. We will instead follow the second approach mentioned above and concentrate on the *semicontinuous model*, as studied in [15], [47], [51], [71], [88], [123], [152]–[154].

2.4. The semicontinuous model. In general, we consider the case when the one-stage cost function $c(\cdot, \cdot)$ is unbounded. Since, for the most part, the criteria considered in this paper are given by a sum of expected costs over the infinite horizon, then, to avoid indeterminate situations, the following conditions will be assumed to hold throughout the paper, unless otherwise indicated.

Assumption 2.1. $c(x, a) \geq 0$ for all $(x, a) \in \mathbf{K}$.

Assumption 2.2. The transition kernel $P(\cdot \mid x, a)$ is *weakly continuous* in (x, a) ; that is, $v(\cdot) \in C_b(\mathbf{S})$ implies that $\int_{\mathbf{S}} v(y)P(dy \mid \cdot, \cdot) \in C_b(\mathbf{K})$.

Assumption 2.3. (i) The multifunction $U(x)$ is upper semicontinuous; (ii) $c(\cdot, \cdot) \in \mathcal{L}(\mathbf{K})$.

Remark 2.1. Concerning Assumption 2.1, note that (for the AC and DC criteria) we must only assume that the cost is bounded below. The assumption that the cost is nonnegative is only made for convenience. Assumption 2.2 is equivalent to $\int v(y)P(dy \mid \cdot, \cdot) \in \mathcal{L}(\mathbf{K})$, for each $v(\cdot) \in \mathcal{L}(\mathbf{S})$ [51, p. 52]. This property is crucial in our development.

Example 2.3. For the nonlinear stochastic system in Example 2.1, assume further that

- (i) \mathbf{A} is compact,
- (ii) For each $x \in \mathbf{S}$, $U(x)$ is closed (and therefore compact), and
- (iii) The system function $F : \mathbf{K} \times \mathbf{W} \rightarrow \mathbf{S}$ is continuous.

If $c(\cdot, \cdot) \in \mathcal{L}(\mathbf{K})$, then, by Remark 2.1, Assumption 2.2 will hold. Furthermore, the assumption on the compactness of \mathbf{A} can be dispensed with if there are compact subsets $\mathbf{K}_1 \subseteq \mathbf{K}_2 \subseteq \dots$ in $\mathbf{S} \times \mathbf{A}$, such that $\mathbf{K} = \bigcup_{n \in \mathbb{N}} \mathbf{K}_n$ and

$$\liminf_{n \rightarrow \infty} \{c(x, a) : (x, a) \in \mathbf{K}_n \setminus \mathbf{K}_{n-1}\} = +\infty,$$

since, in this case, \mathbf{A} can be conveniently compactified; cf. [15, Cor. 8.6.1, p. 210]. Also, the case in which $\mathbf{S} = \mathbb{R}^n$, $\mathbf{A} = \mathbb{R}^m$, and $c(x, a) = x'Qx + a'Ra$, where Q and

R are positive semidefinite and positive definite matrices, respectively, of appropriate dimensions can also be considered by a (one-point) compactification of \mathbf{A} [164, pp. 965–966].

Under Assumptions 2.1–2.3, the *undiscounted dynamic programming map* T given by

$$(2.5) \quad T(v)(x) := \inf_{a \in U(x)} \left\{ c(x, a) + \int_{\mathbf{S}} v(y) P(dy | x, a) \right\} \quad \forall x \in \mathbf{S}$$

maps $\mathcal{L}(\mathbf{S})$ into itself. Also, for $0 < \beta < 1$, the *discounted dynamic programming map* $T_\beta : \mathcal{L}(\mathbf{S}) \rightarrow \mathcal{L}(\mathbf{S})$ is given by

$$(2.6) \quad T_\beta(v) := T(\beta v).$$

The following properties are easily verified.

LEMMA 2.1. *Let $v, v' \in \mathcal{L}(\mathbf{S})$. Then (i) for all $k \in \mathbb{R}$, $T(v + k) = T(v) + k$; (ii) if $v \leq v'$, then $T(v) \leq T(v')$.*

Some key results for the stochastic control problem under a DC criterion are summarized in the following theorem.

THEOREM 2.1. *Under Assumptions 2.1–2.3*

(i) *The following equation, which is called the discounted cost optimality equation (DCOE), holds:*

$$(2.7) \quad J_\beta^*(x) = T_\beta(J_\beta^*)(x) = \inf_{a \in U(x)} \left\{ c(x, a) + \beta \int_{\mathbf{S}} J_\beta^*(y) P(dy | x, a) \right\}, \quad x \in \mathbf{S};$$

(ii) *A policy $\pi^* \in \Pi_{SD}$ is discount optimal if and only if $\pi^*(x)$ attains the infimum in (2.7), for all $x \in \mathbf{S}$;*

(iii) *A discount optimal policy $\pi^* \in \Pi_{SD}$ exists;*

(iv) *Define $T_\beta^0 : \mathcal{L}(\mathbf{S}) \rightarrow \mathcal{L}(\mathbf{S})$ as the identity operator and $T_\beta^k : \mathcal{L}(\mathbf{S}) \rightarrow \mathcal{L}(\mathbf{S})$, $k \in \mathbb{N}$, by $T_\beta^k(f) := T_\beta(T_\beta^{k-1}(f))$. Then, for any $f \in \mathcal{L}_b(\mathbf{S})$,*

$$T_\beta^k(f)(x) \xrightarrow[k \rightarrow \infty]{} J_\beta^*(x) \quad \text{for all } x \in \mathbf{S};$$

(v) *$J_\beta^*(\cdot)$ is nonnegative and lower semicontinuous.*

Remark 2.2. The above results are essentially contained in [15], [51]. The existence of a measurable selector that attains the infimum in (2.7), e.g., the result in (iii) of Theorem 2.1, follows from [15, Prop. 7.33, p. 153], [29], [47, pp. 35–38], [51, §2.6], [88], [139, Thm. 4.1, p. 9], [184, Thm. 9.1, p. 880]. The scheme used in (iv) of Theorem 2.1 to compute $J_\beta^*(\cdot)$ is called the *value iteration* (or successive approximations) algorithm. When the one-stage cost function is bounded, the usual approach is to prove the existence of a unique solution to the DCOE via a contraction mapping theorem [14], [82]. Otherwise, $J_\beta^*(\cdot)$ is not necessarily the only fixed point of T_β ; however, $J_\beta^*(\cdot)$ is the *minimal* fixed point of T_β among the class of nonnegative functions in $\mathcal{L}(\mathbf{S})$ [15, Chap. 5], [173].

3. A sketch of historical development. We now present a brief historical sketch of the development of CMP, with an emphasis on the average cost criterion. The roots of CMP can be traced back to the pioneering work of Wald [186], [187] on sequential analysis and statistical decision functions. In the late 1940s and early

1950s, several investigators formulated the essential concepts of CMP, which are found in their work in sequential game models. A CMP can be viewed as a one-player game. Of particular interest is the work of Bellman and Blackwell [12], Bellman and LaSalle [13], and also Shapley, who formulated the essential mechanism of stochastic dynamic programming and used the theory of contraction mappings [160]. Using his famous heuristic “minimum cost to go,” Bellman showed how powerful the dynamic programming technique was by using it to solve problems in a myriad of settings [9]–[11]. Bellman studied mostly problems with a finite horizon, for which the backward induction approach of dynamic programming suffices to give a complete treatment. The situation is quite different in problems over an infinite horizon. Early work on CMP is also reported in econometrics [4], [49].

Howard [95] was apparently the first to study CMP with an average cost criterion. His *policy iteration* algorithm was the first major computational breakthrough, and his book helped establish CMP as an independent subject of investigation. For CMP with finite state and action spaces, Howard’s policy iteration scheme established the existence of a stationary deterministic policy, optimal in this class only. Derman [38] and Viskov and Shiryaev [183] independently showed that this policy was optimal among all admissible policies. Other computational methods were later proposed. Manne [125] gave a linear programming formulation for the AC criterion, and Wagner [185] later characterized extreme-point optima of the linear program as stationary deterministic policies. White [197] introduced the value iteration (or successive approximations) technique. Excellent accounts of these and other computational methods are given in [14, §5.2] and [137].

On the theoretical side, Blackwell’s seminal paper [18] gave considerable impetus to research in this area, motivating numerous other papers. In [18] Blackwell studied CMP with finite state and action spaces. He considered the DC criterion in great detail and established many important results. In the same paper, he initiated an approach for the AC case, which we will refer to as the *vanishing discount approach*: he treated the AC case as a limit of the DC case, when the discount factor goes to 1, i.e., the discounting effect vanishes. Blackwell established in [18] the existence of a stationary deterministic policy that is discount optimal, for all β sufficiently close to 1. This type of optimality is now called *Blackwell optimality* [14, pp. 336–341]. The relation between the discounted and average case also becomes apparent via Tauberian theorems [87, §4.6]. This fact seems to have been observed first by Gillette [79], who used Tauberian theorems to establish the existence of optimal stationary policies in a stochastic game problem with an AC criterion. Also, using Tauberian theorems, Derman [38] showed that the Blackwell optimal policy found in [18] was also optimal for the AC criterion. Average cost CMP with finite state and arbitrary action spaces were studied under various conditions in the works of [35], [57]–[59], [100].

Blackwell optimal policies do not necessarily exist when the state space is countably infinite [122]. In fact, average optimal policies need not exist in this situation [121], [150]. Similar nonexistence result holds when the state space is finite, but the action space is an arbitrary compact metric space [8]. For such models, the existence of an optimal policy has been proved by Bather [8], Martin-Löf [126], and Feinberg [58], under certain conditions. Derman [39] studied the problem with countable state space, finite action space, and bounded cost. He studied the following equation, which

became known as the *average cost optimality equation* (ACOE):

$$\rho + h(i) = \min_{a \in U(i)} \left\{ c(i, a) + \sum_{j \in \mathcal{S}} P(j | i, a) h(j) \right\},$$

where ρ is a scalar, $h : \mathcal{S} \rightarrow \mathbb{R}$, $\mathcal{S} = \mathbb{N}_0$, and we write $P(j | \cdot, \cdot)$ for $P(\{j\} | \cdot, \cdot)$. He showed that, if the ACOE has a *bounded solution*, i.e., a solution (ρ, h) with $h(\cdot)$ a bounded function, then the stationary deterministic policy realizing the pointwise minimum on the right-hand side of the ACOE is average optimal, and ρ is the minimum average cost. Derman's paper, in conjunction with Derman and Veinott [43], showed that a sufficient condition for the existence of such a solution was that the expected hitting time of a fixed state under *any* stationary deterministic policy is bounded uniformly with respect to the choice of the policy and the initial state. Motivated by Blackwell's work, Taylor [177] extended the vanishing discount approach to obtain a bounded solution for a Markovian sequential replacement problem by studying the asymptotics of the *differential* discounted value function $h_\beta(\cdot) := J_\beta(\cdot) - J_\beta(0)$. Ross [147], [148] refined Taylor's procedure and showed that, under the Derman–Veinott [43] condition, $\{h_\beta(\cdot)\}_{\beta \in (0,1)}$ was uniformly bounded in β . By letting $\beta \uparrow 1$, Ross established that the ACOE had a bounded solution. This made the vanishing discount approach very popular. In subsequent works, many variants of Derman–Veinott recurrence conditions appeared. See [52], [178] for a great variety of such conditions. These conditions are difficult to remove, and counterexamples abound [150]. Actually, it has been shown in [64], in a very general setting, that the uniform boundedness of $\{h_\beta(\cdot)\}_{\beta \in (0,1)}$ in β is also a *necessary* condition for a bounded solution to the ACOE to exist. Cavazos-Cadena [30], [31], under some additional conditions, showed that the existence of bounded solutions to the ACOE necessarily impose a very strong recurrence structure on the model. Lippman [115] studied controlled semi-Markov processes with unbounded cost with both discounted and average cost criteria. Following the vanishing discount approach, he derived results for the average cost case under several restrictive assumptions. Federgruen, Hordijk, and Tijms [53] have explored the same approach.

Hordijk [91] extended many earlier results to countable state space and compact action spaces. He introduced the Lyapunov function method for CMP. He used this method to obtain a (possibly *unbounded*) solution to the ACOE, yielding an optimal policy. However, the Lyapunov function method necessarily imposes a blanket stability of the processes (in the sense of positive recurrence). Such stability is not always met in, e.g., many queueing model applications. In addition, he introduced some new concepts, particularly based on the relation of stochastic dynamic programming with Markov potential theory. There is a vast amount of literature devoted to CMP in several volumes of the Mathematisch Centrum tracts; see [181] and the references therein.

With Hordijk's work, it appeared that a shift away from the vanishing discount approach was necessary. Rosberg, Varaiya, and Walrand [144] treated the average cost criterion as the limiting case of the finite horizon problem, but details of their arguments depend heavily on the specifics of the problem they consider, viz., the control of two queues in tandem with a linear cost structure. Federgruen and Tijms [56] initiated a direct study of the ACOE by a span seminorm method, for bounded costs. This method allows us to obtain useful value iteration algorithms. Later, Federgruen, Schweitzer, and Tijms [55] treated the problem with countable state space

and unbounded costs. Assuming a recurrence condition on the model, they established the existence of a (possibly unbounded) solution to the ACOE, thereby establishing the existence of an optimal stationary deterministic policy.

In a series of papers [20]–[25], Borkar presented a convex analytic approach to treat the problem with countable state space, compact action space, and unbounded cost. This approach can be seen as an extension of the ideas in Manne [125] and Wagner [185]. Borkar stressed the existence of an optimal *stable* stationary deterministic policy, i.e., one that induces a positive recurrent process. While a blanket stability assumption (e.g., of Lyapunov type) may be too restrictive to cover many queueing applications, it nevertheless is desirable that the optimal policy be stable. Borkar showed that, to obtain an optimal stable stationary deterministic policy, either a blanket stability hypothesis or a condition on the cost that penalizes unstable behavior is necessary. He also emphasized the concept of almost sure optimality by a “pathwise” treatment of the problem. A comprehensive account of the convex analytic approach to CMP is given in [26].

After the extensive works of Hordijk, Federgruen et al., and Borkar, it seemed that the vanishing discount approach was not appropriate for many classes of problems with unbounded costs. However, this approach has been revived and generalized to a great extent in [17], [61], [63], [74], [76], [77], [83], [85], [155], [156], [167], [172], [190]. In some of these references, an *inequality* version of the ACOE is studied. In view of the results of [30], [31], and [64], it is clear that a bounded solution to the ACOE is too restrictive, in general. A natural candidate solution is one that is bounded below [28], [76], [85], [155], [156], [172], [190], or one having suitable growth properties [28] or satisfying other conditions [167]. Weber and Stidham [172], [190] studied the problem for queueing systems. Under a penalizing condition on the cost and some structural assumptions, they established the existence of a (possibly unbounded) solution to the ACOE and showed the existence of an optimal stationary deterministic policy. Sennott proceeded along similar lines. She identified very general conditions on the discounted value function so that the vanishing discount approach could successfully be pursued. We refer to [155]–[157], [172], [190] for many interesting examples of queueing systems and to [34] for a comparison of different sets of assumptions. Extensions of these techniques to semi-Markov decisions processes with applications to queueing systems have been reported in [157].

The first attempt to give a description of CMP with more general state and actions spaces was carried out by Karlin [98]. Blackwell [19], Maitra [123], and Strauch [173] studied CMP with a general state space and the discounted cost criterion. Their work was significantly extended by Shreve and Bertsekas in [15], [164], [165]. Feinberg [60] studied CMP with Borel state space and with arbitrary numerical criteria, which include TC, AC, and DC as particular cases. By establishing the convexity of the set of strategic measures (measures of the type \mathcal{P}_μ^π on the canonical space), he established the existence of an ε -optimal $f \in \Pi_{SD}$ for these criteria. De Leve [112]–[114] considered general state and action space CMP in continuous time with an AC criterion, with an emphasis on the ergodic behavior of the processes. Ross [148] used the vanishing discount approach to study CMP with an AC criterion, general state space, finite action space, and bounded cost function. He showed that, if the family of differential discounted value functions $\{h_\beta(\cdot)\}_{\beta \in (0,1)}$ is equicontinuous and uniformly bounded, then the ACOE admits a bounded solution, yielding an optimal stationary deterministic policy. Ross also introduced the concept of minorant. He showed that,

if there exists a state $x_0 \in \mathbf{S}$ and $\alpha > 0$ such that

$$P(x_0 | x, a) > \alpha \quad \text{for all } a \in U(x), \quad x \in \mathbf{S},$$

then the average cost case could be reduced to a discounted one. This was greatly extended in the work of Gubenko and Statland [80] (see also [43]). They showed that, under similar minorant (or majorant) conditions, a contraction map, with respect to the sup norm, could be defined on $\mathcal{M}_b(\mathbf{S})$, which would yield a bounded solution to the ACOE. They also obtained bounded solutions to the ACOE under continuity and boundedness conditions, which guarantee that $\{h_{\beta_n}(\cdot)\}$, with $\beta_n \uparrow 1$, is uniformly bounded and equicontinuous; thus a similar approach as in [148] can be followed. Georgin [72], [73] also explored this approach, under some ergodicity conditions. Tijms [179] and Hübner [96] directly studied the ACOE, under some ergodicity assumptions, by showing that the undiscounted dynamic programming map is a contraction on $\mathcal{M}_b(\mathbf{S})$, with respect to the span seminorm. For an excellent presentation of these methods and the type of ergodicity conditions used, see [82, §3.3]. Wijngaard [201], [202] and Kumar [101] studied the problem under Doeblin's condition using an operator theoretic method. Under several conditions, Kurano [104] obtained solutions to the ACOE and also showed the existence of an average optimal stationary deterministic policy. Also, in [105]–[107], he obtained the existence of an optimal stationary deterministic policy under Doeblin's condition. For a comprehensive presentation of the different recurrence conditions used for the above purposes, see [86].

The study of *partially observable controlled Markov processes* (POCMP) was initiated independently by various authors [5], [46], [50], [161], [162]. The reduction to models with complete information (see §7) was exhibited for various cases in [5], [138], [151], [205]. The study of finite state space POCMP with an AC criterion was initiated by Sondik [170]. Transforming the problem into an equivalent *completely observable* (CO) problem with Borel state space, Sondik tried to cast the problem in the framework of Ross [148] but did not show equicontinuity of $\{h_\beta(\cdot)\}_{\beta \in (0,1)}$. Ross [150], Wang [189], and White [191] showed this equicontinuity condition for specific scalar replacement problems. Ohnishi, Mine, and Kawai [132] studied a multistate replacement problem by using concavity properties of $h_\beta(\cdot)$. Platzman studied the general problem of finite state and action space POCMP, also by using concavity properties of the functions $h_\beta(\cdot)$. Under certain reachability conditions, he proved that the family $\{h_\beta(\cdot)\}_{\beta \in (0,1)}$ is uniformly bounded. However, even though this family may not be equicontinuous with respect to the Euclidean metric, he showed that it is equi-Lipschitzian with respect to some other appropriate metric, thus putting the problem within the framework of Ross [148]. Fernández-Gaucherand, Arapostathis, and Marcus [62], [63] followed a different approach to the problem, using the concepts of invariant sets of a CMP and controlled sub-Markov processes. This approach allows us to consider POCMP with countable state and observation spaces. Borkar [26] also studied the problem via his convex analytic approach.

4. Finite state space. In this section, we will consider models with a finite state space. Initially, we restrict our attention to the case when \mathbf{A} is a finite set; models with compact action space will be discussed at the end of the section.

4.1. Finite action spaces. Let $\mathbf{S} = \{1, \dots, k\}$. In this case, Π_{SD} is finite. This fact plays a crucial role in the analysis for the average cost problem. For a policy $\pi \in \Pi$, let $J_\beta(\pi)$ denote the vector $(J_\beta(1, \pi), \dots, J_\beta(k, \pi))^T$; similarly, we define $J_N(\pi)$, $J(\pi)$,

J_β^* , J^* , and J_N^* . For a stationary deterministic policy $f \in \Pi_{SD}$, let $P(f)$ denote the transition matrix of the corresponding process and

$$c(f) := (c(1, f(1)), \dots, c(k, f(k)))^T.$$

Also, the (i, j) th entry in the n th power of the transition matrix $P(f)$ will be denoted by $P_{ij}^n(f)$ or $P^n(f)(i, j)$. It is well known that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} P^n(f) := P^*(f)$$

exists [18], [87, Chap. 4], [137], where $P^0(f) = I$ (the $k \times k$ identity matrix). Using the theory of stochastic matrices, the following results can be proved. For details, see [8], [18], [87], [137].

THEOREM 4.1. For each $f \in \Pi_{SD}$,

- (i) $J(f) = P^*(f)c(f)$;
- (ii) The number of linearly independent equations in $(I - P(f))w = c(f) - J(f)$ is k minus the number of communicating classes in $P(f)$;
- (iii) The equations

$$(4.1) \quad (I - P(f))w = c(f) - v,$$

$$(4.2) \quad P^*(f)w = 0$$

have solutions $v = J(f)$ and $w = w(f)$, where

$$w(f) := (I - P(f) + P^*(f))^{-1}(I - P^*(f))c(f);$$

- (iv) $v = J(f)$ and $w = w(f)$ are the unique solutions to (4.1) and (4.2) for which $v(s) = v(s')$ if s and s' are in the same communicating class of $P(f)$, and $v(s) = J(s, f)$ if state s is transient in $P(f)$.

Remark 4.1. (a) It is easily seen from the above theorem that if, under an $f \in \Pi_{SD}$, the process is irreducible or unichain (see [87]), then $J(\cdot, f)$ is constant.

- (b) The matrix

$$H(f) := (I - P(f) + P^*(f))^{-1}(I - P^*(f))$$

is called the *deviation matrix*. It plays a fundamental role in the analysis. For the discounted case, $J_\beta(f) = (I - \beta P(f))^{-1}c(f)$. Analogous results can be developed for the average cost case using $H(f)$. The following result, due to Miller and Veinott [127] and Lamond and Puterman [110], can be proved using the spectral theory of stochastic matrices.

THEOREM 4.2. Let $\beta \in [0, 1)$ and $\lambda = (1 - \beta)\beta^{-1}$. Let $f \in \Pi_{SD}$ and ν be the eigenvalue of $P(f)$ less than one with largest modulus. If $0 \leq \lambda \leq 1 - |\nu|$, then

$$(4.3) \quad (\lambda I + I - P)^{-1} = \lambda^{-1}P^*(f) + \sum_{n=0}^{\infty} (-\lambda)^n H^{n+1}(f)$$

and

$$(4.4) \quad J_\beta(f) = (1 + \lambda) \left[\lambda^{-1}P^*(f)c(f) + \sum_{n=0}^{\infty} (-\lambda)^n H^{n+1}(f)c(f) \right].$$

Remark 4.2. (a) The quantity $h(f) := H(f)c(f)$ plays a crucial role in the analysis of the problem. It is called the *bias* or *transient cost*. It can be easily seen from the Neumann series expansion of $(I - P(f) + P^*(f))^{-1}$ [137] that, for $s \in \mathbf{S}$,

$$h(f)(s) = E_s^f \left[\sum_{t=0}^{\infty} (c(X_t, f(X_t)) - J(X_t, f)) \right].$$

From the above representation, $h(f)$ can be interpreted as the expected total cost for a CMP with cost $c - J$. If $P(f)$ is aperiodic, the distribution of X_t converges to a limiting distribution, so eventually $c(X_t, f(X_t))$ and $J(X_t, f)$ will differ very little. Thus, $h(f)$ can be thought of as the expected total cost “until convergence” or the expected total cost during the “transient” phase of the evolution of the process [137].

(b) Howard [95] has shown that

$$J_N(f) = NJ(f) + h(f) + o(1).$$

Therefore, as N becomes large, for each $s \in \mathbf{S}$, $J_N(f)$ approaches a straight line with slope $J(f)$ and intercept $h(f)$. When $J(f)(s)$ is constant, $J_N(s) - J_N(s')$ approaches $h(f)(s) - h(f)(s')$, so that $h(f)$ is the asymptotic relative difference of starting the process in two states s and s' . That is why $h(f)$ is often referred to as the relative value. See [14, pp. 304–308], [36] for a good discussion of these matters.

(c) Expansion (4.4) extends Blackwell’s expansion [18].

(d) Using expansion (4.4), the following important result is immediate.

COROLLARY 4.1. For $f \in \Pi_{SD}$, $J(f) = \lim_{\beta \uparrow 1} (1 - \beta)J_\beta(f)$.

Following Blackwell [18] and Derman [40], we now prove the following existence results.

THEOREM 4.3. *There exists an $f \in \Pi_{SD}$ that is discount optimal for all β sufficiently close to 1 and is also optimal for the average cost criterion.*

Proof. For each $f \in \Pi_{SD}$ and $s \in \mathbf{S}$, $J_\beta(s, f)$ is obviously an analytic function of β . Let $\{\beta_n\}$, $0 < \beta_n < 1$ be a sequence such that $\beta_n \uparrow 1$. For a fixed n , let $f_n \in \Pi_{SD}$ be β_n -discount optimal (see Theorem 2.1). Since Π_{SD} is a finite set, the sequence $\{f_n\}$ must contain at least one $f^* \in \Pi_{SD}$ that occurs infinitely often. Let $\{\beta_{n_k}\}$ be a subsequence of $\{\beta_n\}$ such that $\beta_{n_k} \uparrow 1$ and $f^* = f_{n_1} = f_{n_2} = \dots$. Then, for every $g \in \Pi$, $J_{\beta_{n_k}}(f^*) \leq J_{\beta_{n_k}}(g)$. Since all coordinates of $J_\beta(f^*)$ and $J_\beta(g)$ are analytic functions of β , it follows that

$$J_\beta(f^*) \leq J_\beta(g)$$

for all β near 1. Since this holds for all $g \in \Pi$, it follows that f^* is β -discount optimal for all β near 1. We next show that f^* is average optimal. Let $\pi \in \Pi$. Then

$$(1 - \beta_{n_k})J_{\beta_{n_k}}(f^*) \leq (1 - \beta_{n_k})J_{\beta_{n_k}}(\pi), \quad k = 1, 2, \dots$$

Therefore, letting $k \rightarrow \infty$ and using Theorem 4.1 and a standard Tauberian theorem (Theorem A.2 in the Appendix), it follows that

$$\begin{aligned} J(f^*) &= \lim_{\beta \uparrow 1} (1 - \beta)J_\beta(f^*) \\ &= \lim_{k \rightarrow \infty} (1 - \beta_{n_k})J_{\beta_{n_k}}(f^*) \\ &\leq \limsup_{k \rightarrow \infty} (1 - \beta_{n_k})J_{\beta_{n_k}}(\pi) \leq J(\pi), \end{aligned}$$

and the proof is complete. \square

We now briefly mention three numerical approaches. For details, we refer to [14], [137], [180], among others. Our presentation follows [137].

Value iteration. We assume that, under any $f \in \Pi_{SR}$, the corresponding chain is unichain and aperiodic. For any positive integer N , the finite horizon value function J_N^* satisfies the equation

$$(4.5) \quad J_{N+1}^* = \min_{f \in \Pi_{SD}} \{c(f) + P(f)J_N^*\}.$$

Equation (4.5) can act as an iteration equation with $J_0^* \equiv 0$ as the initial condition. Let $f_{N+1}^* \in \Pi_{SD}$ realize the minimum in (4.5). We can treat $(1/N)J_N^*$ and f_N^* as our guesses for J^* and an average optimal policy. Then $J_N^* - NJ^*$ converges as $N \rightarrow \infty$. Also, there exists an integer N_0 such that, for any $N \geq N_0$, any $f \in \Pi_{SD}$ that attains the minimum in (4.5) is average optimal. However, this property does not yield an error estimate and hence fails to provide a stopping rule for the iteration scheme. To this end, with $h = (h(1), \dots, h(k))$, we let

$$L(h) := \min_{x \in \mathbf{S}} \{Th(x) - h(x)\}, \quad U(h) := \max_{x \in \mathbf{S}} \{Th(x) - h(x)\}.$$

It can be shown that [137]

$$\min_{x \in \mathbf{S}} \{J_N^*(x) - J_{N-1}^*(x)\} \leq J^* \leq \max_{x \in \mathbf{S}} \{J_N^*(x) - J_{N-1}^*(x)\}$$

and

$$L(J_{N-1}^*) \leq L(J_N^*) \leq J^* \leq U(J_N^*) \leq U(J_{N-1}^*).$$

Furthermore, $\lim_{N \rightarrow \infty} \{U(J_N^*) - L(J_N^*)\} = 0$. Thus, an average ε -optimal policy can be found by stopping the value iteration when

$$U(J_N^*) - L(J_N^*) < \varepsilon.$$

There are other variants of this approach; see [54], [56], and [96].

Linear programming. To simplify our presentation, we will assume that, under any $f \in \Pi_{SR}$, the corresponding process is irreducible. Let $P(f)$ denote the transition matrix of the process, and $\eta(f) \in \mathcal{P}(\mathbf{S})$ its invariant measure. Then, for any $s \in \mathbf{S}$, $J(s, f) = J(f)$, a constant, and

$$J(f) = \sum_{s \in \mathbf{S}} \sum_{a \in U(s)} c(s, a) f(s, a) \eta(f)(s).$$

Therefore, the average cost problem can be reduced to the following linear programming problem:

$$(4.6a) \quad \text{minimize} \quad \sum_{s \in \mathbf{S}} \sum_{a \in U(s)} c(s, a) x(s, a)$$

subject to

$$(4.6b) \quad x(s, a) \geq 0, \quad s \in \mathbf{S}, \quad a \in U(s),$$

$$(4.6c) \quad \sum_{s \in \mathbf{S}} \sum_{a \in U(s)} x(s, a) = 1,$$

$$(4.6d) \quad \sum_{a \in U(s)} x(s, a) = \sum_{s' \in \mathbf{S}} \sum_{a \in U(s')} x(s', a) P(s' | s, a).$$

Under the irreducibility assumption, the simplex method can be employed to find an optimal stationary deterministic policy. This formulation is due to Manne [125].

Policy improvement. We work under the irreducibility assumption. The dual to the linear program (4.6a)–(4.6d) is the problem of finding variables g and $h(s)$, $s \in \mathcal{S}$, to

$$(4.7a) \quad \text{maximize } g$$

subject to

$$(4.7b) \quad g + \sum_{s' \in \mathcal{S}} (\delta(s, s') - P(s' | s, a)) h(s) \leq c(s, a),$$

$(s, a) \in \mathcal{S} \times U(s)$, where $\delta(s, s')$ is the Kronecker delta.

The functional equation

$$(4.8) \quad g + h(s) = \min_{a \in U(s)} \left\{ c(s, a) + \sum_{s' \in \mathcal{S}} P(s' | s, a) h(s) \right\}$$

is equivalent to (4.7a), (4.7b) under the irreducibility assumption and is the average cost optimality equation [87]. We will discuss this equation in detail in the next section. It will be shown that an $f \in \Pi_{SD}$ is optimal if and only if f realizes the pointwise minimum in (4.8), and then g is the optimal average cost. This suggests the following iteration algorithm.

- (i) Let $n = 1$. Choose $f_n \in \Pi_{SD}$. Let $h_n(s) \equiv 0$ for all $s \in \mathcal{S}$.
- (ii) Find a solution g_n and $h_n(s)$ of the following equation:

$$g_n + h_n(s) = c(s, f_n(s)) + \sum_{s' \in \mathcal{S}} P(s' | s, f_n(s)) h_n(s).$$

- (iii) For each $s \in \mathcal{S}$, compute

$$\phi_n(s) = \min_{a \in U(s) \setminus \{f_n(s)\}} \left\{ c(s, a) + \sum_{s' \in \mathcal{S}} P(s' | s, a) h_n(s') \right\} - g_n - h_n(s).$$

If $\phi_n(s) \geq 0$ for all $s \in \mathcal{S}$, then f_n is average optimal and g_n is the optimal average cost. If $\phi_n(s) < 0$ for some $s \in \mathcal{S}$, then pick $a \in U(s)$ such that

$$c(s, a) + \sum_{s' \in \mathcal{S}} P(s' | s, a) h_n(s') - g_n - h_n(s) < 0.$$

Define $f_{n+1} \in \Pi_{SD}$ as $f_{n+1}(s) = a$ and $f_{n+1}(\cdot) = f_n(\cdot)$, otherwise. Then f_{n+1} yields a lower average cost. Since Π_{SD} is finite, the policy improvement scheme converges in a finite number of steps.

4.2. Compact action spaces. We now consider the problem where the action set \mathbf{A} is not finite but a compact metric space. In this situation, an optimal policy may not exist; see [51, p. 178, Ex. 1]. Note that here Π_{SD} is no longer finite. Under certain ergodicity assumptions, Martin-Löf [126] and Feinberg [57] have proved the existence of an optimal $f \in \Pi_{SD}$. We will discuss various ergodicity assumptions on a countable state space in detail in the next section. First, we focus on ε -optimal policies established by Chitashvili [35] and Feinberg [58]; see [51, Chap. 7].

THEOREM 4.4. *Under Assumptions 2.1–2.3, for every $\varepsilon > 0$, there exists an ε -optimal $f \in \Pi_{SD}$.*

Proof (Sketch). For $f \in \Pi_{SD}$, let $J(f)$ be as in Theorem 4.1. For $i \in \mathbf{S}$, let

$$(4.9) \quad \tilde{J}(i) = \inf_{f \in \Pi_{SD}} J(f)(i).$$

Clearly, $J^*(i) \leq \tilde{J}(i)$, for each $i \in \mathbf{S}$. Corresponding to $i \in \mathbf{S}$, select an $f_i \in \Pi_{SD}$ such that

$$(4.10) \quad J(f_i)(i) \leq \tilde{J}(i) + \varepsilon.$$

The set $\tilde{\mathbf{A}} = \{f_i(j) : i, j \in \mathbf{S}\}$ is obviously finite. Taking the action set to be $\tilde{\mathbf{A}}$, the preceding results can be applied to the finite CMP $(\mathbf{S}, \tilde{\mathbf{A}}, P, c)$. For this model, there exists a stationary deterministic policy, say f^* , which is average optimal. Thus

$$(4.11) \quad J(f^*)(i) \leq J(f_i)(i) \leq \tilde{J}(i) + \varepsilon \quad \text{for each } i \in \mathbf{S}.$$

Let

$$\rho^*(i) := \limsup_{\beta \uparrow 1} (1 - \beta)J_\beta^*(i).$$

Then, by Theorem A.2, in the Appendix,

$$(4.12) \quad \rho^*(i) \leq J^*(i) \quad \text{for each } i \in \mathbf{S}.$$

By (4.11), it suffices to show that $J^*(i) = \tilde{J}(i)$, for each $i \in \mathbf{S}$. From (4.12), it then suffices to show that $\rho^*(i) \geq \tilde{J}(i)$. For each $\beta \in (0, 1)$, let $f_\beta \in \Pi_{SD}$ be β -discount optimal. Let f be a limit point of f_β as $\beta \uparrow 1$. Then using (4.4) (which is valid in this case as well) and Assumptions 2.1–2.3, it can be shown that

$$\rho^*(i) \geq J(f)(i) \geq \tilde{J}(i). \quad \square$$

Concerning the existence of an optimal policy, we state the following result.

THEOREM 4.5. *Let Assumptions 2.1 and 2.2 hold and further assume that $c(\cdot)$ is continuous on Π_{SR} and that, under any $f \in \Pi_{SR}$, the corresponding chain is unichain. Then there exists an optimal policy in Π_{SD} .*

The result is almost immediate from the fact that, under the unichain assumption, $P^*(\cdot)$, and therefore also $J(\cdot)$, is continuous on Π_{SR} [91, Lemma 10.2]. For further details, including the convergence of a policy improvement algorithm, see [94].

5. Countable state space. The average cost problem becomes much more complicated when the state space is countable. Maitra [121] has given a counterexample that shows that there need not exist an optimal policy. In [122] Maitra has studied a particular problem in which there does not exist any policy that is β -discount optimal for all β sufficiently close to 1. Flynn [68] has constructed a more dramatic counterexample. In his example, there exists an average optimal policy in Π_{SD} . Nevertheless he exhibits an $f \in \Pi_{SD}$ and a $\beta_0 \in (0, 1)$ such that f is β -discount optimal for all $\beta \in (\beta_0, 1)$, but it is not average optimal. Fisher and Ross [67] have presented a counterexample that shows that the optimal policy need not be stationary or deterministic. We refer to [150] for several other counterexamples. It is apparent that the average

cost problem is closely related to the ergodic behavior of the process, and it is well known that the ergodic theory of Markov processes on a countable state space is much more involved than on a finite state space; for example, a Markov process on a finite state space cannot be null recurrent. Another vital difference in this case is that the number of stationary deterministic policies is no longer finite. To study the ergodic theory, some recurrence conditions are necessary. There are many such conditions available in the literature [26], [52], [178]; we will survey a few representative ones.

In what follows, the state space $\mathbf{S} = \{0, 1, 2, \dots\}$. For each $i \in \mathbf{S}$, the action space $U(i)$ is a prescribed compact metric space. We will always assume that, for fixed $i, j \in \mathbf{S}$, $c(i, \cdot)$, $P(i | j, \cdot)$, are continuous. These conditions can be weakened or dropped in several places, as will be clear from the specific context.

Derman [38] studied the ACOE that, with ρ a scalar and $h : \mathbf{S} \rightarrow \mathbb{R}$, takes the following form:

$$(5.1) \quad \rho + h(i) = \min_{a \in U(i)} \left\{ c(i, a) + \sum_{j \in \mathbf{S}} P(j | i, a) h(j) \right\}.$$

A solution to (5.1) is a pair (ρ, h) satisfying it.

Suppose that $f \in \Pi_{SD}$ is a minimizing selector in (5.1). Then (5.1) becomes

$$(5.1') \quad \rho + h(i) = c(i, f(i)) + \sum_{j \in \mathbf{S}} P(j | i, f(i)) h(j).$$

Equation (5.1') asserts that, apart from ρ , the cost if the process stops now equals the expected cost if it continues under the policy f for just one more period. We can give a similar interpretation to (5.1). Hence, we may think that ρ is the average cost under f and that no other $f \in \Pi_{SD}$ has a smaller average cost. Thus, the function h in (5.1) is roughly a measure of how much we are prepared to pay to stop the process, though continuing to pay an average cost ρ in the future [141] (cf. Remark 4.2(a)). Therefore, the function h may be viewed as a cost potential. Also, by a stochastic representation of h , using (5.1) and (5.1'), h is indeed a potential. Hordijk [91] has pursued this line of thought in great detail, which we will discuss later.

We start with a characterization of optimal policies.

THEOREM 5.1. *If the ACOE has a solution (ρ, h) satisfying*

$$(5.2) \quad \lim_{t \rightarrow \infty} \frac{1}{t} E_i^\pi h(X_t) = 0 \quad \forall \pi \in \Pi_{SD}, \quad \forall i \in \mathbf{S},$$

then there exists an $f \in \Pi_{SD}$ such that

$$\rho = J(i, f) = J^*(i) \quad \forall i \in \mathbf{S}.$$

Moreover, an $f \in \Pi_{SD}$ is average optimal if, for each $i \in \mathbf{S}$,

$$(5.3) \quad c(i, f(i)) + \sum_{j \in \mathbf{S}} P(j | i, f(i)) h(j) = \min_{a \in U(i)} \left\{ c(i, a) + \sum_{j \in \mathbf{S}} P(j | i, a) h(j) \right\},$$

and, conversely, if an $f \in \Pi_{SD}$ is average optimal and the corresponding chain is irreducible and positive recurrent, then (5.3) holds.

Proof. Let $f \in \Pi_{SD}$ satisfy (5.3). Then, since

$$E_i^f [h(X_{t+1}) | \mathfrak{F}_t] = \sum_{j \in \mathbf{S}} P(j | X_t, f(X_t)) h(j),$$

it follows from (5.1) and (5.3) that

$$(5.4) \quad \rho + h(X_t) = c(X_t, f(X_t)) + E_i^f[h(X_{t+1}) \mid \mathfrak{F}_t].$$

Summing (5.4) from $t = 0$ to $N - 1$, dividing by N , and taking expectations, we obtain

$$\rho = \frac{1}{N} E_i^f \left[\sum_{t=0}^{N-1} c(X_t, f(X_t)) \right] + \frac{E_i^f[h(X_N)] - h(i)}{N}.$$

Next, letting $N \rightarrow \infty$ and using (5.2) yields

$$\rho = \lim_{N \rightarrow \infty} \frac{1}{N} E_i^f \left[\sum_{t=0}^{N-1} c(X_t, f(X_t)) \right].$$

On the other hand, if π is any other policy, we can show using the same arguments that

$$\rho \leq \limsup_{N \rightarrow \infty} \frac{1}{N} E_i^\pi \left[\sum_{t=0}^{N-1} c(X_t, A_t) \right].$$

Hence, f is average optimal. Conversely, let $f \in \Pi_{SD}$ be average optimal and suppose that the corresponding chain is irreducible and positive recurrent. If f does not satisfy (5.3), then there exist $i_0 \in \mathcal{S}$, $a_0 \in U(i_0)$ and $\delta > 0$ such that

$$(5.5) \quad \begin{aligned} c(i_0, f(i_0)) + \sum_{j \in \mathcal{S}} P(j \mid i_0, f(i_0)) h(j) \\ = c(i_0, a_0) + \sum_{j \in \mathcal{S}} P(j \mid i_0, a_0) h(j) + \delta. \end{aligned}$$

Let $f' \in \Pi_{SD}$ be defined as follows:

$$f'(i) = \begin{cases} f(i) & \text{if } i \neq i_0, \\ a_0 & \text{if } i = i_0. \end{cases}$$

Then, using (5.5) along with irreducibility and positive recurrence, it is easily seen that $J(i_0, f') < J(i_0, f)$, which contradicts the average optimality of f . \square

Remark 5.1. (a) We say that (5.1) admits a bounded solution if $h(\cdot)$ is bounded. If the ACOE has a bounded solution, then (5.2) is clearly satisfied; moreover, using the martingale stability theorem [117, p. 53], it can be shown that the $f \in \Pi_{SD}$ selecting the minimum in (5.3) is sample path average optimal [72].

(b) Various extensions of last assertion of Theorem 5.1 have been obtained by Sennott [158].

Derman and Veinott [43] have prescribed a certain recurrence condition that ensures that (5.1) admits a bounded solution. We will discuss it later in this section. The ACOE resembles the dynamic programming equation, and Theorem 5.1 is analogous to a dynamic programming characterization of an optimal policy. However, the dynamic programming heuristic does not lead directly to the ACOE. Taylor [177] developed a vanishing discount approach for a particular problem, which was extended for the general case by Ross [147]–[150]. Our presentation here follows Ross [150]. As noted earlier, the average case can in some sense be treated as the limiting case of

the discounted problem as the discount factor approaches 1. The discounted value function $J_\beta^*(\cdot)$ satisfies the DCOE (cf. Theorem 2.1)

$$J_\beta^*(i) = \min_{a \in U(i)} \left\{ c(i, a) + \beta \sum_{j \in \mathbf{S}} P(j | i, a) J_\beta^*(j) \right\},$$

and a β -discounted optimal policy selects a minimizing action. One possible way of finding an average optimal policy might be to choose the actions minimizing

$$\lim_{\beta \rightarrow 1} \left\{ c(i, a) + \beta \sum_{j \in \mathbf{S}} P(j | i, a) J_\beta^*(j) \right\}.$$

However, this limit need not exist and indeed would often be infinite for all actions. The situation can nevertheless be salvaged by considering a “differential” discounted value function, i.e., $h_\beta(i) := J_\beta^*(i) - J_\beta^*(0)$, where $0 \in \mathbf{S}$ is an arbitrary, fixed state. The function $h_\beta(\cdot)$ satisfies

$$(5.6) \quad (1 - \beta)J_\beta^*(0) + h_\beta(i) = \min_{a \in U(i)} \left\{ c(i, a) + \beta \sum_{j \in \mathbf{S}} P(j | i, a) h_\beta(j) \right\}.$$

From (5.6) it is now apparent that (5.1) can be derived under certain conditions by letting $\beta \rightarrow 1$. We state here a simple result [150], despite the fact that it also holds under weaker hypotheses (see Theorem 5.9).

THEOREM 5.2. *Suppose that there exists a constant $K > 0$ such that $|h_\beta(i)| \leq K$, for all $\beta \in (0, 1)$ and $i \in \mathbf{S}$. Then*

- (i) *The ACOE admits a bounded solution (ρ, h) ;*
- (ii) *For some sequence $\beta_n \rightarrow 1$, $h(i) = \lim_{n \rightarrow \infty} h_{\beta_n}(i)$, $i \in \mathbf{S}$;*
- (iii) *$\lim_{\beta \rightarrow 1} (1 - \beta)J_\beta^*(i) = \rho$ for any $i \in \mathbf{S}$.*

Proof. Let $\beta_n \uparrow 1$ be given. By the uniform boundedness of $h_\beta(\cdot)$, using a diagonalization procedure, we can find a subsequence, which for simplicity we also denote by β_n , such that $h_{\beta_n}(i) \rightarrow h(i)$ for each $i \in \mathbf{S}$, where $h(\cdot)$ is a bounded function. Again, since $(1 - \beta_n)J_{\beta_n}^*(0)$ is bounded, there is a further subsequence $\beta_{n_k} \uparrow 1$ such that

$$\lim_{k \rightarrow \infty} (1 - \beta_{n_k})J_{\beta_{n_k}}^*(0)$$

exists. Part (i) of the theorem then follows from (5.6) and an application of the dominated convergence theorem. Furthermore, by Theorem 5.1, ρ is the minimum average cost. Since the above results are independent of the sequence chosen, (iii) then follows. \square

Remark 5.2. It has been shown [64] that, if the ACOE has a bounded solution, then there exists a constant $K > 0$ such that $|h_\beta(i)| \leq K$ for all $\beta \in (0, 1)$, $i \in \mathbf{S}$.

5.1. Bounded costs. In this section, we assume that $c(\cdot, \cdot)$ is bounded. Ross [150] has proved that under a Derman–Veinott [43] type recurrence condition (see (5.7), below), the uniform boundedness hypothesis of Theorem 5.2 is satisfied.

THEOREM 5.3. *Let $f \in \Pi_{SD}$ and let $\{X_t\}$ be the corresponding state process. Let*

$$\tau = \min\{t \geq 1 : X_t = 0\}.$$

If there exists a $K > 0$ such that

$$(5.7) \quad E_i^f[\tau] < K$$

for all $f \in \Pi_{SD}$ and all $i \in \mathbf{S}$, then $h_\beta(i)$ is bounded uniformly in $\beta \in (0, 1)$ and $i \in \mathbf{S}$.

Proof. Let $\beta \in (0, 1)$ and $f_\beta \in \Pi_{SD}$ be β -discount optimal. We have

$$(5.8) \quad \begin{aligned} J_\beta^*(i) &= E_i^{f_\beta} \left[\sum_{t=0}^{\infty} \beta^t c(X_t, f_\beta(X_t)) \right] \\ &= E_i^{f_\beta} \left[\sum_{t=0}^{\tau-1} \beta^t c(X_t, f_\beta(X_t)) \right] + E_i^{f_\beta} \left[\sum_{t=\tau}^{\infty} \beta^t c(X_t, f_\beta(X_t)) \right] \\ &\leq M E_i^{f_\beta}[\tau] + J_\beta^*(0) E_i^{f_\beta}[\beta^\tau], \end{aligned}$$

where M is a bound on $c(\cdot, \cdot)$. From (5.7) and (5.8), it follows that

$$(5.9) \quad J_\beta^*(i) - \beta J_\beta^*(0) \leq MK.$$

Also, from (5.8) and applying Jensen's inequality, we obtain

$$J_\beta^*(i) \geq J_\beta^*(0) E_i^{f_\beta}[\beta^\tau] \geq J_\beta^*(0) \beta^K.$$

Therefore,

$$(5.10) \quad \begin{aligned} J_\beta^*(0) - J_\beta^*(i) &\leq (1 - \beta^K) J_\beta^*(0) \\ &\leq (1 - \beta^K) \frac{M}{1 - \beta} \leq MK. \end{aligned}$$

The desired result follows from (5.9) and (5.10). \square

After the work of Derman [38], Derman and Veinott [43], and Ross [147], [148], several recurrence conditions have appeared [178]. We explore a few representative ones.

Let $f \in \Pi_{SD}$. For a finite set $A \subset \mathbf{S}$, let

$$(5.11) \quad \tau_A = \min\{t \geq 1 : X_t \in A\}.$$

Assumption 5.1. There is a finite $A \subset \mathbf{S}$ and a constant $K > 0$ such that $E_i^f[\tau_A] < K$ for all $i \in \mathbf{S}$ and $f \in \Pi_{SD}$. Furthermore, for any $f \in \Pi_{SD}$ the corresponding process does not have two disjoint invariant sets.

Assumption 5.2. There exists a constant $K > 0$, and, for every $f \in \Pi_{SD}$, there is a state $j(f) \in \mathbf{S}$ such that

$$E_i^f[\tau_{\{j(f)\}}] < K \quad \forall i \in \mathbf{S}.$$

Assumption 5.3 (simultaneous Doeblin). There is a finite set \mathbf{A} , an integer $n \geq 1$ and a scalar $\alpha > 0$ such that

$$\sum_{j \in \mathbf{A}} P(j \mid i, f(i)) \geq \alpha$$

for all $i \in \mathbf{S}$ and all $f \in \Pi_{SD}$. Furthermore, for any $f \in \Pi_{SD}$, the corresponding process does not have two disjoint invariant sets.

Assumption 5.4 (scrambling). There is an integer $n \geq 1$ and a scalar $\alpha > 0$ such that, for any $f \in \Pi_{SD}$,

$$\sum_{j \in \mathbf{S}} \min\{P_{i_1,j}^n(f), P_{i_2,j}^n(f)\} \geq \alpha \quad \forall i_1, i_2 \in \mathbf{S}.$$

Assumption 5.5 (ergodicity). There is an integer $n \geq 1$ and a scalar $\rho > 0$ such that, for each $f \in \Pi_{SD}$, there exists an $\eta(f) \in \mathcal{P}(\mathbf{S})$ for which

$$\sum_j |P_{ij}^m(f) - \eta(f)(j)| \leq 2(1 - \rho)^{\lfloor m/n \rfloor}$$

for all $i \in \mathbf{S}$ and $m \geq 1$, where $\lfloor x \rfloor$ denotes the largest integer not exceeding x .

Remark 5.3. Clearly Assumptions 5.1 and 5.2 are generalizations of the Derman–Veinott condition. Hordijk [91] has proved the existence of a bounded solution to the ACOE using Assumption 5.1. Under Assumption 5.5, for each $f \in \Pi_{SD}$, $\eta(f)$ is the unique invariant measure of the corresponding process.

Federgruen, Hordijk, Tijms [52] have established the following theorem.

THEOREM 5.4. *Assumptions 5.1–5.3 are equivalent. Also, if for any $f \in \Pi_{SD}$ the corresponding process is aperiodic, then Assumptions 5.1–5.5 are equivalent.*

Remark 5.4. Under any one of Assumptions 5.1–5.5, Federgruen, Hordijk, and Tijms [52] have established the existence of a bounded solution to the ACOE by extending the vanishing discount approach of Taylor and Ross.

We have thus far seen several recurrence conditions which are sufficient for the ACOE to admit a bounded solution. Cavazos-Cadena [30], [31] has dealt with the converse question of what are the necessary recurrence conditions for the ACOE to have a bounded solution. He has obtained the following result. Consider the following assumption.

Assumption 5.6. There exists a constant $K > 0$ such that, for each bounded and measurable $c : \mathbf{S} \times \mathbf{A} \rightarrow \mathbb{R}$ and every collection $\{U(i) : i \in \mathbf{S}\}$, $U(i) \subset \mathbf{A}$, there exist $\rho \in \mathbb{R}$ and $h : \mathbf{S} \rightarrow \mathbb{R}$ bounded that solve (5.1) and satisfy $\|h\| \leq K\|c\|$, where $\|\cdot\|$ is the sup norm.

THEOREM 5.5. *Assumptions 5.2 and 5.6 are equivalent.*

The proof follows by an application of the uniform boundedness principle. For details and other variants, we refer to [30], [31]. Thus, an assumption on the existence of a bounded solution to the ACOE necessarily imposes a strong recurrence structure on the system. Also, note that Assumption 5.6 involves not just one CMP but a family of CMP (one for each c and $\{U(i)\}$). Since it is equivalent to Assumptions 5.1–5.3 and under aperiodicity conditions to Assumptions 5.1–5.5, it follows that Assumptions 5.1–5.5 are too strong for many important applications. In fact, there are interesting situations [20] in which these conditions are not satisfied, but for which we can find average optimal stationary deterministic policies.

Ross [148] has proved that, under the following recurrence condition, the AC can be reduced to an appropriate DC. Therefore, in view of Theorem 2.1, the problem can be resolved in this case.

THEOREM 5.6. *If there exists a constant $\alpha > 0$ such that*

$$P(0 \mid i, a) \geq \alpha > 0$$

for all $i \in \mathbf{S}$, $a \in U(i)$, then the AC can be reduced to an appropriate DC.

Proof. Let

$$\tilde{P}(j | i, \cdot) = \begin{cases} (1 - \alpha)^{-1}P(j | i, \cdot) & \text{for } j \neq 0, \\ (1 - \alpha)^{-1}(P(0 | i, \cdot) - \alpha) & \text{for } j = 0. \end{cases}$$

Let $\tilde{J}_\beta^*(\cdot)$ denote the β -discounted value function for the CMP with cost $c(\cdot, \cdot)$ and transition law $\tilde{P}(\cdot | \cdot, \cdot)$. Then it is easily verified that, for each $i \in \mathbf{S}$,

$$\alpha \tilde{J}_{1-\alpha}^*(0) + \tilde{J}_{1-\alpha}^*(i) = \min_{a \in U(i)} \left\{ c(i, a) + \sum_{j \in \mathbf{S}} P(j | i, a) \tilde{J}_{1-\alpha}^*(j) \right\}.$$

Let $f \in \Pi_{SD}$ be $(1 - \alpha)$ -discount optimal for the modified CMP. It follows from Theorem 5.1 that f is AC-optimal for the original CMP, and the optimal average cost is $\alpha \tilde{J}_{1-\alpha}^*(0)$. \square

Remark 5.5. Note that, if the ACOE has a bounded solution (ρ, h) , then ρ is the optimal average cost for any initial condition. Hence, the existence of a bounded solution to the ACOE suggests that some kind of “unchainedness” is in effect, since, for the multichain case, the average cost would, in general, depend on the initial condition. The multichain version of the ACOE is

$$(5.12a) \quad \min_{a \in U(i)} \sum_{j \in \mathbf{S}} P(j | i, a) \rho(j) = \rho(i),$$

$$(5.12b) \quad \rho(i) + h(i) = \min_{a \in U_1(i)} \left\{ c(i, a) + \sum_{j \in \mathbf{S}} P(j | i, a) h(j) \right\},$$

where

$$(5.12c) \quad U_1(i) = \left\{ a \in U(i) : \min_{a \in U(i)} \sum_{j \in \mathbf{S}} P(j | i, a) \rho(j) = \rho(i) \right\}.$$

This equation has been studied by Zijm [208] for countable state space. For more general state spaces, it was extensively studied much earlier by Yushkevich [204] (see also [51]); this work will be discussed in the next section.

If (5.12) has a bounded solution $\rho(i), h(i)$, where both ρ and h are bounded functions, then we can show, as before, that $\rho(i)$ is the optimal average cost starting from state $i \in \mathbf{S}$ and a minimizing selector in (5.12) yields an average optimal stationary deterministic policy. Under a certain “geometric convergence condition,” Zijm [208] has established the existence of a bounded solution to (5.12). Under the additional assumptions that under any stationary deterministic policy the corresponding process has at most a finite number of ergodic classes, he has shown that the geometric convergence condition is equivalent to a number of recurrence conditions of the type Assumptions 5.1–5.5.

Hordijk [91] establishes the existence of an average optimal $f \in \Pi_{SD}$ without utilizing the ACOE. Let Π_{SD} be endowed with the product topology. Then Π_{SD} is compact and metrizable. Let us consider the following assumptions.

Assumption 5.7. For each $f \in \Pi_{SD}$ and $i \in \mathcal{S}$, there exists a measure $\eta_i(f) \in \mathcal{P}(\mathcal{S})$ such that $\eta_i(f)(j) = \lim_{N \rightarrow \infty} (1/N) \sum_{n=0}^{N-1} P^n(f)(i, j)$.

Assumption 5.8. $f \mapsto \eta_i(f)$ is continuous for any $i \in \mathcal{S}$.

Assumption 5.9. For each $i \in \mathcal{S}$, $\{\eta_i(f) : f \in \Pi_{SD}\}$ is tight (for a definition of tightness, see [134, Def. 3.1, p. 28]).

Assumption 5.10. For each $f \in \Pi_{SD}$, the corresponding process is recurrent.

Assumption 5.11. For each $f \in \Pi_{SD}$, the corresponding process does not have disjoint-invariant sets.

Assumption 5.12. $\{P(f)(i, \cdot) : i \in \mathcal{S}, f \in \Pi_{SD}\}$ is tight.

It is easy to see that Assumptions 5.7 and 5.8 imply that, for each $i \in \mathcal{S}$, $\{\eta_i(f) : f \in \Pi_{SD}\}$ is compact. Hence, in particular, Assumptions 5.7 and 5.8 imply Assumption 5.9. By definition, Assumption 5.9 implies Assumption 5.7. Also, it can easily be shown that Assumptions 5.9 and 5.11 imply Assumption 5.8, and that Assumption 5.12 implies 5.9. However, Assumption 5.12 may be easier to verify.

THEOREM 5.7. *Each of the following five combinations of assumptions is sufficient for the existence of an average optimal $f \in \Pi_{SD}$: (Assumption 5.7, Assumption 5.8), (Assumption 5.9, Assumption 5.10), (Assumption 5.9, Assumption 5.11), (Assumption 5.10, Assumption 5.12), (Assumption 5.11, Assumption 5.12).*

Remark 5.6. The main idea behind the proof of this theorem can be traced back to the proof of Theorem 4.3. We give the main points and skip the details. Let $\beta_n \in (0, 1)$ be a sequence such that $\beta_n \uparrow 1$, let $f_{\beta_n} \in \Pi_{SD}$ be β_n -discount optimal, and f_∞ be a limit point of $\{f_{\beta_n}\}$ in Π_{SD} . Suppose that $\rho^*(i)$ is a scalar satisfying $(1 - \beta_n)J_{\beta_n}^*(i) \rightarrow \rho^*(i)$, for each $i \in \mathcal{S}$ (along a suitable subsequence). Then, by using Tauberian and ergodic theorems, we deduce that $J^*(i) = \rho^*(i)$ and f_∞ is average optimal under (Assumption 5.7, Assumption 5.8). Under (Assumption 5.9, Assumption 5.10), f_∞ is average optimal for initial states $i \in \tilde{\mathcal{S}} := \bigcup_i \text{supp}(\eta_i(f_\infty))$, where ‘‘supp’’ denotes the support. Then by Assumption 5.10 there exists an \tilde{f} such that the corresponding process starting from any $i \in \mathcal{S} \setminus \tilde{\mathcal{S}}$ reaches $\tilde{\mathcal{S}}$. Set

$$\tilde{f}(i) = \begin{cases} \tilde{f}(i) & \text{if } i \notin \tilde{\mathcal{S}}, \\ f_\infty(i) & \text{if } i \in \tilde{\mathcal{S}}. \end{cases}$$

It follows that \tilde{f} is average optimal. The other cases can be dealt with in a similar manner.

5.2. Unbounded costs. We have thus far considered bounded costs only. There are practical situations (e.g., in queueing systems) where the cost is typically unbounded. We assume that $c \geq 0$ (cf. Assumption 2.1). Let us now consider the ACOE for unbounded c . Note that the boundedness of c , did not play any role in the proof of Theorem 5.1. For unbounded c , the ACOE is unlikely to admit a bounded solution.

Lippman [115], [116] has studied controlled semi-Markov processes with unbounded costs. He has placed polynomial bounds on the movement of the process in one transition. He has made a further assumption that there exists an $f \in \Pi_{SD}$ such that both the mean first passage times and mean first passage costs from any state i to state zero under the policy are finite. Moreover, if $f \in \Pi_{SD}$ is close to β -discount optimal for a sequence of discount factors, then it is AC-optimal. Lippman has employed the vanishing discount approach of Taylor and Ross to establish the existence of a solution (ρ, h) to the ACOE with h satisfying (5.2), thereby establishing the existence of an

average optimal $f \in \Pi_{SD}$. He has also given some examples from queueing systems where his conditions are satisfied. However, his condition on the β -discounted value function appears to be very difficult to verify.

Hordijk [91] has used a Lyapunov stability condition to establish the existence of an average optimal $f \in \Pi_{SD}$.

Assumption 5.13 (Lyapunov condition). Let

$$\tilde{P}(f)(i, j) = \begin{cases} P(f)(i, j), & j \neq 0, \\ 0, & j = 0. \end{cases}$$

There exists a function $w : \mathcal{S} \rightarrow \mathbb{R}_+$ such that, for all $i \in \mathcal{S}$,

- (i) $c(i, f(i)) + 1 + \sum_j \tilde{P}(f)(i, j)w(j) \leq w(i)$, for all $f \in \Pi_{SD}$;
- (ii) $\sum_j P(f)(i, j)w(j)$ is continuous in f ;
- (iii) $\lim_{n \rightarrow \infty} \sum_j \tilde{P}^n(f)(i, j)w(j) = 0$.

THEOREM 5.8. *Under the above Lyapunov condition, there exists an AC-optimal $f \in \Pi_{SD}$.*

Proof (Sketch). Let $f \in \Pi_{SD}$. For $i \in \mathcal{S}$, we define $\tau_i = \min\{t \geq 1 : X_t = i\}$, where X_t is governed by f . Then, under Assumption 5.13, using the standard techniques of stochastic Lyapunov function method [91], [108], the following results can be proved:

$$(5.13) \quad E_i^f[\tau_0] \leq w(i),$$

$$(5.14) \quad E_i^f \left[\sum_{t=0}^{\tau_0-1} c(X_t, f(X_t)) \right] \leq w(i).$$

Indeed, with $n \in \mathbb{N}$ and $n > 1$,

$$\begin{aligned} E_i^f[w(X_{n \wedge \tau_0}) \mid \mathfrak{F}_{n \wedge \tau_0}] - w(i) &= -E_i^f \left[\sum_{t=0}^{n \wedge \tau_0 - 1} E_i^f[w(X_{t+1}) \mid X_t] - w(X_t) \right] \\ &\leq -E_i^f[n \wedge \tau_0], \end{aligned}$$

where the last inequality is due to Assumption 5.13. Hence, $E_i^f[n \wedge \tau_0] \leq w(i)$, and, letting $n \uparrow \infty$, (5.13) follows. Also, (5.14) can be proved along the same lines. By an ergodic theorem [133],

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{1}{N} E_0^f \left[\sum_{t=0}^{N-1} c(X_t, f(X_t)) \right] &= (E_0^f \tau_0)^{-1} E_0^f \left[\sum_{t=0}^{\tau_0-1} c(X_t, f(X_t)) \right] \\ &=: \rho(f). \end{aligned}$$

Let $\rho^* := \inf_{f \in \Pi_{SD}} \rho(f)$. Then $\rho^* \leq w(0)$. Define

$$h(i) = \inf_{f \in \Pi_{SD}} E_i^f \left[\sum_{t=0}^{\tau_0-1} (c(X_t, f(X_t)) - \rho^*) \right].$$

Then $h(0) = 0$. Using (5.13), (5.14), and Assumption 5.13(iii), it can be shown that (ρ^*, h) is a solution of the ACOE with h satisfying (5.2), and the desired result follows. \square

Remark 5.7. (a) Note that by Assumption 5.13(i) the cost function c does not grow faster than the Lyapunov function w . Thus, there is a restriction on the growth of c imposed by w . In CMP, $w(i) = i$, $w(i) = i^2$ are typical examples of Lyapunov functions [91]. In the latter case, for example, we can treat only those unbounded cost functions that do not grow faster than quadratic functions.

(b) Assumption 5.13(iii) is crucial in showing that the cost potential h satisfies $\lim_{t \rightarrow \infty} (1/t) E_i^f h(X_t) = 0$, for all $f \in \Pi_{SD}$, and $i \in \mathcal{S}$.

Federgruen, Hordijk, and Tijms [53] have extended Hordijk's results by replacing the single attracting point $\{0\}$ by a finite set $K \subset \mathcal{S}$. Their main assumption is the following: There exists a finite set $K \subset \mathcal{S}$ such that, for each initial state $i \in \mathcal{S}$, the suprema over the mean hitting time of K and mean hitting costs are finite. This, in turn, is equivalent to the existence of a Lyapunov function $w : \mathcal{S} \rightarrow \mathbb{R}_+$ satisfying Assumption 5.13(i), where now \tilde{P} is defined as

$$\tilde{P}(f)(i, j) = \begin{cases} P(f)(i, j), & j \notin K, f \in \Pi_{SD}, \\ 0, & j \in K. \end{cases}$$

Under the additional assumptions that Assumption 5.13(ii) and (iii) hold, and the ‘‘communication condition’’ that for any $f \in \Pi_{SD}$ the corresponding process has no two disjoint invariant sets, they have established the existence of a solution (ρ, h) to the ACOE by employing the vanishing discount approach and have shown that h satisfies (5.2). This work has been further extended by Federgruen, Schweitzer, and Tijms [55]. They have dropped the unchainedness assumption in [53]. Instead, they assume that any state can be reached from any other state via some policy. Under this and other conditions in [53], they have established the existence of a solution (ρ, h) to the ACOE, with h satisfying (5.2). They have deviated from the vanishing discount approach and have, instead, utilized Tychonoff's fixed point theorem in their analysis. We again note that, in all these investigations, a restrictive growth condition on the cost function is imposed, as noted in Remark 5.7.

The Lyapunov stability condition necessarily imposes a blanket stability (i.e., positive recurrence) of certain states (cf. (5.13)), which may be very restrictive. On the other hand, (5.2) is not easy to verify in general and, indeed, may not hold in the case of many queueing models [141]. Another generalization of the boundedness of the solution of the ACOE could be boundedness from below. This will be the case if the cost function has some ‘‘monotone’’ properties, which naturally arise in various queueing models. This line of thought has been pursued in various ways in [24], [28], [74], [76], [77], [141], [142], [155], [156], [172], [190].

Sennott [155], [156] has prescribed very general conditions in this direction. We will now briefly describe them. Consider the following assumptions.

Assumption 5.14. For every $i \in \mathcal{S}$ and every $\beta \in (0, 1)$, $J_\beta^*(i) < \infty$.

Assumption 5.15. There exists a nonnegative integer L such that

$$h_\beta(i) := J_\beta^*(i) - J_\beta^*(0) \geq -L.$$

Assumption 5.16. There exists a function $M : \mathcal{S} \rightarrow \mathbb{R}_+$ such that $h_\beta(i) \leq M(i)$ for all $i \in \mathcal{S}$ and any $\beta \in (0, 1)$. For every $i \in \mathcal{S}$, there exists an $a(i) \in U(i)$ such that

$$\sum_j P(j \mid i, a(i)) M(j) < \infty.$$

THEOREM 5.9. *Under Assumptions 5.14–5.16, there exists an AC-optimal $f \in \Pi_{SD}$.*

Proof. Let $\beta_n \in (0, 1)$ be such that $\beta_n \uparrow 1$. Let f_{β_n} be β_n -discount optimal. Let f be a limit point of f_{β_n} as $n \rightarrow \infty$. To simplify the notation, all subsequences of β_n will also be denoted by β_n . By Assumption 5.16 and a diagonal argument, there exists a function $h : \mathbf{S} \rightarrow \mathbb{R}$ such that $\lim_{n \rightarrow \infty} h_{\beta_n}(\cdot) = h(\cdot)$. By Assumption 5.15, $h(\cdot) \geq -L$. Let $\rho : \mathbf{S} \rightarrow \mathbb{R}_+$ be a function such that $\lim_{n \rightarrow \infty} (1 - \beta_n)J_{\beta_n}^*(i) = \rho(i)$. Using Assumption 5.16, it is easy to see that $\rho(i) = \rho^*$, a constant. Now, for $i \in \mathbf{S}$,

$$(5.15) \quad (1 - \beta_n)J_{\beta_n}^*(0) + h_{\beta_n}(i) = c(i, f_{\beta_n}(i)) + \beta_n \sum_{j \in \mathbf{S}} P(j | i, f_{\beta_n}(i))h_{\beta_n}(j).$$

Fix an $i \in \mathbf{S}$. Add L to both sides to make $(h_{\beta_n}(i) + L) \geq 0$ and take “lim inf” on both sides of (5.15). Then, by Fatou’s lemma and the assumption of continuity of $P(j | i, \cdot)$, we conclude that

$$\rho^* + h(i) \geq c(i, f(i)) + \sum_j P(j | i, f(i))h(j).$$

Since $h(\cdot)$ is bounded below, the proof of Theorem 5.1 can be modified to show that $J(i, f) \leq \rho^*$. By Theorem A.2 in the Appendix, $J(i, \pi) \geq \rho^*$ for any $\pi \in \Pi$. Hence, $J(i, f) = J^*(i) = \rho^*$, and f is AC-optimal. \square

Remark 5.8. (a) From the above proof, it is clear that if ρ is a scalar, $h : \mathbf{S} \rightarrow \mathbb{R}$ is bounded below, and

$$(5.16) \quad \rho + h(i) \geq \min_{a \in U(i)} \left\{ c(i, a) + \sum_j P(j | i, a)h(j) \right\},$$

then ρ is the optimal average cost, and any $f \in \Pi_{SD}$ selecting the minimum on the right-hand side of (5.16) is AC-optimal. In this case, we may replace the ACOE by an *average cost optimality inequality* (ACOI), viz., (5.16).

(b) If, for each $i \in \mathbf{S}$, $U(i)$ is finite, then, in the above proof, $f_{\beta_n}(i) = f(i)$ for large n . Then we can write, for large n ,

$$\rho + h(i) = c(i, f(i)) + \beta_n \sum_j P(j | i, f(i))h_{\beta_n}(j).$$

By Fatou’s lemma,

$$\rho + h(i) \geq c(i, f(i)) + \sum_j P(j | i, f(i))h(j).$$

Consider the stronger assumption, below.

Assumption 5.17. Assumption 5.16 holds, and $\sum_j P(j | i, a)M(j) < \infty$, for all $a \in \mathbf{A}$ and $i \in \mathbf{S}$.

Under Assumption 5.17, using dominated convergence, it is easy to see that

$$\rho + h(i) = \min_{a \in U(i)} \left\{ c(i, a) + \sum_j P(j | i, a)h(j) \right\},$$

and we obtain the ACOE. If, for each $i \in \mathbf{S}$, there is a finite set $R_i \subset \mathbf{S}$ such that $P(j | i, \cdot) = 0$ for $j \notin R_i$, then Assumption 5.17 will obviously hold. Such a condition is satisfied for systems whose dynamics have a nearest-neighbour motion property [28].

(c) If there exists an $f \in \Pi_{SD}$, under which the process is ergodic, irreducible with an invariant measure $\eta(f) \in \mathcal{P}(\mathbf{S})$, and $\sum_i c(i, f(i))\eta(f)(i) < \infty$, then Assumptions 5.14 and 5.16 hold. Assumption 5.15 holds if $J_\beta^*(i)$ is increasing in i . Direct conditions implying Assumptions 5.14–5.17 can be found in [28], [32], [34], [76], [77], [155], [156], [172], [190]. See also [141], [142].

(d) Let $f \in \Pi_{SD}$ be a policy that attains the minimum on the right-hand side of (5.16). Fix an $i \in \mathbf{S}$. If the chain under f is positive recurrent at i , then we can show that equality holds at i in (5.16). However, the lack of positive recurrence at i may lead to strict inequality in (5.16). Cavazos-Cadena [33] had exhibited an example to demonstrate this. He has further shown in his example [33] that Assumptions 5.14–5.16 are satisfied, but the ACOE does not admit any solution.

5.3. The convex analytic approach. We will now describe Borkar’s convex analytic approach for the average cost case [20]–[26]. The convex analytic approach to the AC-problem is a natural extension of the linear programming approach when the state/action spaces are no longer finite. In this approach, we view the control problem as the problem of minimizing a linear functional on the convex set of “ergodic occupation measures,” to be defined shortly [20]–[26]. This approach can also be used to treat other standard cost criteria, but it may be more involved for treating cases such as the DC criterion. On the other hand, it is more flexible and powerful for certain other purposes, e.g., pathwise average cost, constrained optimization problem, among others. Since the techniques involved here are entirely different from what we have thus far followed, we will embark on a more detailed discussion.

By replacing each $U(i)$ with $\prod_k U(k)$ and $P(j | i, \cdot)$ by its composition with the projection $\prod_k U(k) \rightarrow U(i)$, we may and will assume that the $U(i)$ ’s are replicas of a fixed compact metric space \mathbf{A} . We say that an $f \in \Pi_{SR}$ is *stable* if the corresponding process is positive recurrent. We will assume that, under an $f \in \Pi_{SR}$, the process has \mathbf{S} as its single communicating class. (This can be relaxed in some cases; see [26] for a discussion on this.) Therefore, f will have a unique invariant measure $\eta(f) \in \mathcal{P}(\mathbf{S})$ satisfying

$$\eta(f)P(f) = \eta(f).$$

Let Π_{SSR} denote the space of stable stationary policies. Π_{SSD} is defined analogously. For an $f \in \Pi_{SSR}$, denote by $\hat{\eta}(f) \in \mathcal{P}(\mathbf{S} \times \mathbf{A})$ the “ergodic occupation measure” defined by

$$\int_{\mathbf{S} \times \mathbf{A}} g d\hat{\eta}(f) = \sum_{i \in \mathbf{S}} \eta(f)(i) \int_{\mathbf{A}} g(i, a) f(i)(da)$$

for $g \in C_b(\mathbf{S} \times \mathbf{A})$. We will consider the sample path average cost optimality, which is stronger than the usual AC-optimality. Let

$$I_R = \{\hat{\eta}(f) : f \in \Pi_{SSR}\}, \quad I_D = \{\hat{\eta}(f) : f \in \Pi_{SSD}\}.$$

Note that $\hat{\eta}(f)$ can only be defined for an $f \in \Pi_{SSR}$. To consider optimality in Π , we will need to consider the following empirical processes. Let $\pi \in \Pi$ and let (X_t, A_t) be the corresponding processes with initial law $\mu \in \mathcal{P}(\mathbf{S})$. Define the $\mathcal{P}(\mathbf{S} \times \mathbf{A})$ -valued empirical process $\{\nu_t\}_{t \geq 1}$ by

$$(5.17) \quad \nu_t(C \times D) = \frac{1}{t} \sum_{s=0}^{t-1} I\{X_s \in C, A_s \in D\}, \quad t \geq 1,$$

for C, D Borel in \mathbf{S}, \mathbf{A} , respectively. Let $\bar{\mathbf{S}} = \mathbf{S} \cup \{\infty\}$ be the one-point compactification of \mathbf{S} . By abuse of notation, we may identify ν_t with the element of $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ that restricts to it on $\mathbf{S} \times \mathbf{A}$. Since $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ is compact, $\{\nu_t\}$, viewed as a sequence of $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ -valued random variables, converges to a sample path dependent compact limit set in $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$. We characterize this set in Lemma 5.1, below, the statement of which calls for some new notation. Note that any element $\nu \in \mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ can be decomposed as

$$(5.18) \quad \nu(B) = \delta_\nu \nu'(B \cap (\mathbf{S} \times \mathbf{A})) + (1 - \delta_\nu) \nu''(B \cap (\{\infty\} \times \mathbf{A}))$$

for B Borel in $\bar{\mathbf{S}} \times \mathbf{A}$, $\delta_\nu \in [0, 1]$ is uniquely specified and $\nu' \in \mathcal{P}(\mathbf{S} \times \mathbf{A})$ (respectively, $\nu'' \in \mathcal{P}(\{\infty\} \times \mathbf{A})$) is uniquely specified if $\delta_\nu > 0$ (respectively, $\delta_\nu < 1$). We may render ν', ν'' unique at all times by imposing an arbitrary fixed choice thereof when $\delta_\nu = 0$, respectively, 1.

LEMMA 5.1. *Outside a set of zero probability (with respect to \mathcal{P}_μ^π), the following holds: For any limit point ν of $\{\nu_t\}$ in $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ for which $\delta_\nu > 0$,*

$$\nu' = \hat{\eta}(f)$$

for some $f \in \Pi_{SSR}$.

Proof. By the martingale stability theorem [117, p. 53],

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=1}^t [I\{X_s = i\} - E_\mu^\pi [I\{X_s = i\} \mid \mathfrak{F}_{s-1}]] \\ = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=1}^t \left[I\{X_s = i\} - \sum_{j \in \mathbf{S}} P(i \mid j, A_{s-1}) I\{X_{s-1} = j\} \right] \\ = \lim_{t \rightarrow \infty} \left[\nu_t(\{i\} \times \mathbf{A}) - \int P(i \mid \cdot, \cdot) d\nu_t \right] \\ = 0 \quad \text{a.s.,} \end{aligned}$$

for each $i \in \mathbf{S}$. Consider a sample path outside the set of zero probability on which the above fails for any $i \in \mathbf{S}$. Then, for any ν as in the statement of the lemma, we must have

$$\nu'(\{i\} \times \mathbf{A}) \geq \int P(i \mid \cdot, \cdot) d\nu', \quad i \in \mathbf{S}.$$

Note that an inequality is obtained here, since the second term on the right-hand side of (5.18) is obviously nonnegative. Summing over $i \in \mathbf{S}$ on both sides, it follows that equality must hold. Decomposing ν' as $\nu'(i, da) = \bar{\nu}(i) f(i)(da)$, where $\bar{\nu} \in \mathcal{P}(\mathbf{S})$ is the marginal on \mathbf{S} and $i \mapsto f(i) \in \mathcal{P}(\mathbf{A})$ is a version of the regular conditional law that defines an element of Π_{SR} , we obtain

$$\bar{\nu}(i) = \sum_{j \in \mathbf{S}} \bar{\nu}(j) P(f)(i, j).$$

Hence, $\bar{\nu} = \eta(f)$, and the conclusion follows. \square

LEMMA 5.2. *The sets I_R and I_D are closed; also, I_R is convex and has its extreme points in I_D .*

Proof. Let $\hat{\eta}(f_n) \in I_R$ and $\hat{\eta}(f_n) \rightarrow \nu$ for some ν in $\mathcal{P}(\mathbf{S} \times \mathbf{A})$. Then, for all $i \in \mathbf{S}$,

$$\hat{\eta}(f_n)(\{i\} \times \mathbf{A}) = \int P(i | \cdot, \cdot) d\hat{\eta}(f_n), \quad n \geq 1.$$

Letting $n \rightarrow \infty$, $\nu(\{i\} \times \mathbf{A}) = \int P(i | \cdot, \cdot) d\nu$. Now argue as in the proof of the preceding lemma to conclude that $\nu = \hat{\eta}(f)$ for some $f \in \Pi_{SSR}$. This proves that I_R is closed. The proof that I_D is closed is similar. Let $f_1, f_2 \in \Pi_{SSR}$ and $0 \leq \lambda \leq 1$. Define $f \in \Pi_{SSR}$ as follows:

$$f(i) = \frac{\lambda\eta(f_1)(i)f_1(i) + (1-\lambda)\eta(f_2)(i)f_2(i)}{\lambda\eta(f_1)(i) + (1-\lambda)\eta(f_2)(i)}.$$

Then using the properties of invariant measures, it is not difficult to see that

$$\begin{aligned} \eta(f) &= \lambda\eta(f_1) + (1-\lambda)\eta(f_2), \\ \hat{\eta}(f) &= \lambda\hat{\eta}(f_1) + (1-\lambda)\hat{\eta}(f_2), \end{aligned}$$

showing that I_R is convex. Now let $f \in \Pi_{SSR}$ be such that, for some $i_0 \in \mathbf{S}$ and $0 < \lambda < 1$, there exist $\phi_1, \phi_2 \in \mathcal{P}(\mathbf{A})$ such that

$$\begin{aligned} \int P(\cdot | i_0, a) f(i_0)(da) &= \lambda \int P(\cdot | i_0, a) \phi_1(da) + (1-\lambda) \int P(\cdot | i_0, a) \phi_2(da), \\ \int P(\cdot | i_0, a) \phi_1(da) &\neq \int P(\cdot | i_0, a) \phi_2(da). \end{aligned}$$

Define $f_1, f_2 \in \Pi_{SSR}$ as

$$f_i(j) = \begin{cases} f(j), & j \neq i_0, \\ \phi_i, & j = i_0. \end{cases}$$

Then it can be shown [24] that $f_1, f_2 \in \Pi_{SSR}$, and any two of $\eta(f)$, $\eta(f_1)$, $\eta(f_2)$ are distinct from each other. Let $b \in (0, 1)$ be such that

$$\lambda = b\eta(f_1)(i_0) / (b\eta(f_1)(i_0) + (1-b)\eta(f_2)(i_0)).$$

Then we can argue as before to conclude that $\hat{\eta}(f) = b\hat{\eta}(f_1) + (1-b)\hat{\eta}(f_2)$. Therefore $\hat{\eta}(f)$ is not an extreme point of I_R . This implies that, for $\hat{\eta}(f')$ to be an extreme point of I_R , $P(\cdot | i, a)$ must be constant over $a \in \text{supp}(f'(i))$, for each $i \in \mathbf{S}$. Hence, $P(f'') = P(f')$, for all $f'' \in \Pi_{SSR}$ such that $\text{supp}(f''(i)) \subset \text{supp}(f'(i))$, for each $i \in \mathbf{S}$. In this case, $\eta(f'') = \eta(f')$. Suppose that for some i , say $i = 1$, there exist $\alpha \in (0, 1)$ and $\phi'_1, \phi'_2 \in \mathcal{P}(\mathbf{A})$, $\phi'_1 \neq \phi'_2$, such that $f'(1) = \alpha\phi'_1 + (1-\alpha)\phi'_2$. Define $f'_1, f'_2 \in \Pi_{SSR}$ by

$$f'_k = \begin{cases} \phi'_k & \text{if } i = 1, \\ f'(i) & \text{if } i \neq 1, \end{cases} \quad k = 1, 2.$$

It follows that $\eta(f') = \eta(f'_1) = \eta(f'_2)$. It is also easy to check that

$$\begin{aligned} \hat{\eta}(f') &= \alpha\hat{\eta}(f'_1) + (1-\alpha)\hat{\eta}(f'_2), \\ \hat{\eta}(f'_1) &\neq \hat{\eta}(f'_2), \end{aligned}$$

which contradicts the extremality of $\hat{\eta}(f')$. Hence, $f'(1)$ must be a Dirac measure. Applying this argument to each $i \in \mathbf{S}$, we deduce that $f \in \Pi_{SSD}$. From this, it follows that the extreme points of I_R lie in I_D . \square

We now proceed to show the existence of a sample path average cost optimal $f \in \Pi_{SSD}$. It is clear that a blanket stability condition or some condition on the cost that penalizes unstable behavior is required to give the desired existence. For example, consider the case where $c(i, a) = \exp(-i)$, which rewards unstable behavior. Clearly, the cost for any $f \in \Pi_{SSR}$ is almost surely positive. On the other hand, provided that $\Pi_{SSR} \neq \Pi_{SR}$, there exists an unstable policy in Π_{SR} that results in an almost-sure zero cost and is, therefore, optimal (the hypothesis that under some $f \in \Pi_{SR}$ the process has \mathbf{S} as its single communicating class plays a crucial role in this assertion). We want to rule out this possibility, as stability is a very desirable property of a policy. We wish to find conditions under which our goal will be achieved. Let $f \in \Pi_{SSR}$. Define

$$\rho(f) := \int c d\hat{\eta}(f), \quad \rho^* := \inf_{f \in \Pi_{SSR}} \rho(f).$$

Note that, under $f \in \Pi_{SSR}$, $J(i, f) = \rho(f)$ for each $i \in \mathbf{S}$. We consider two sets of hypotheses.

Assumption 5.18 (the near-monotonicity condition). It holds that

$$\liminf_{i \rightarrow \infty} \min_{a \in \mathbf{A}} c(i, a) > \rho^*.$$

Intuitively, Assumption 5.18 penalizes the drift of the process away from some finite set, requiring the optimal policy to exert some kind of a “centripetal force” pushing the process back toward this finite set. Thus, the optimal policy gains the desired stability property. If $c(i, a) = k(i)$ for some $k : \mathbf{S} \rightarrow \mathbb{R}_+$ and $k(i)$ is increasing, then this condition will automatically be satisfied. Such penalizing conditions quite often occur in queueing applications (see [20], [155], [156], [172], [190]).

Assumption 5.19 (stability condition (cf. Assumptions 5.7–5.12)). $\Pi_{SR} = \Pi_{SSR}$ and I_R is compact.

Assumption 5.19'. Equivalent conditions to Assumption 5.19 are

- (i) $\Pi_{SD} = \Pi_{SSD}$ and I_D is compact;
- (ii) The mean return times to a prescribed state (say 0) are uniformly integrable over all $f \in \Pi_{SR}$;
- (iii) This is the same as (ii), but with Π_{SD} replacing Π_{SR} .

THEOREM 5.10. *Under Assumption 5.18 or Assumption 5.19, there exists an $f \in \Pi_{SSD}$, which is sample path average cost optimal in Π_{SR} .*

Proof. From Lemma 5.2, it can be shown by an application of Choquet’s theorem [25], [26] that, if $\nu \mapsto \int c d\nu$ attains its minimum on I_R , it will do so for an $f \in \Pi_{SD}$. Under Assumption 5.19, it can be shown that $f \mapsto \hat{\eta}(f)$ is continuous. Therefore, the desired result follows under Assumption 5.19. We next consider the case under Assumption 5.18. Let $f_n \in \Pi_{SR}$ be such that $\rho(f_n) \downarrow \rho^*$. By identifying $\hat{\eta}(f_n)$ with the element of $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ that restricts to it on $\mathbf{S} \times \mathbf{A}$ for each n and then dropping to a subsequence if necessary, we may assume that $\hat{\eta}(f_n) \rightarrow \nu$ in $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ for some ν . Let $n \rightarrow \infty$ in the equation

$$\hat{\eta}(f_n)(\{j\} \times \mathbf{A}) = \int P(j | \cdot, \cdot) d\hat{\eta}(f_n), \quad j \in \mathbf{S}$$

and argue as in Lemma 5.1 to conclude that, for ν' as in (5.18), $\delta_\nu > 0$ implies that

$$\nu'(\{j\} \times \mathbf{A}) = \int P(j | \cdot, \cdot) d\nu', \quad j \in \mathbf{S}.$$

Decomposing ν' as $\nu'(i, da) = \bar{\nu}(i)f(i)(da)$, $i \in \mathbf{S}$, we have $\bar{\nu} = \eta(f)$ and therefore $\nu' = \hat{\eta}(f)$. Let $c_m = c \wedge m$ for $m \geq 1$ and pick $\varepsilon > 0$ such that Assumption 5.18 continues to hold with $\rho^* + \varepsilon$ in place of ρ^* . Then

$$\begin{aligned} \rho^* &= \lim_{n \rightarrow \infty} \int c d\hat{\eta}(f_n) \\ &\geq \lim_{n \rightarrow \infty} \int c_m d\hat{\eta}(f_n) \\ &\geq \delta_\nu \int c_m d\hat{\eta}(f) + (1 - \delta_\nu)((\rho^* + \varepsilon) \wedge m). \end{aligned}$$

Letting $m \rightarrow \infty$,

$$\rho^* \geq \delta_\nu \rho^* + (1 - \delta_\nu)(\rho^* + \varepsilon).$$

This is possible only if $\delta_\nu = 1$ and $\int c d\hat{\eta}(f) = \rho^*$. \square

The above theorem, however, does not ensure optimality of the cost-minimizing policy in I_R with respect to arbitrary policies. For the near-monotone case, this can be resolved without any further assumptions, but, for the stable case, we need the following.

Assumption 5.20. If $\tau = \min\{t \geq 1 : X_t = 0\}$, then

$$\sup_{\pi \in \Pi} E_0^\pi[\tau^2] < \infty.$$

Remark 5.9. Assumption 5.20 clearly implies Assumption 5.19. The converse need not be true, as can be shown by an explicit example [24]. Some sufficient conditions for Assumption 5.20 are (i) a Lyapunov condition [28], which we will describe shortly (cf. Theorem 5.11), (ii) the strong uniform recurrence condition of Doeblin and its variants [178], and (iii) the condition that there exist an $N < \infty$ for which

$$\sup_{\pi \in \Pi} \sup_i \mathcal{P}_i^\pi(\tau \geq N) < 1,$$

where τ is as above.

THEOREM 5.11. *Under Assumption 5.18 or Assumption 5.20, there exists an $f \in \Pi_{SD}$, which is sample path average cost optimal.*

Proof. Under Assumption 5.20, it can be shown [26] that the processes ν_t as defined in (5.17) are tight over Π . Therefore, δ_ν as in the statement of Lemma 5.1 may be taken to be 1. This resolves the case under Assumption 5.20. Under (A5.18), let ν be a limit point of $\{\nu_t\}$ in $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ along some subsequence. Then, as in the proof of Theorem 5.9, it can be shown that

$$(5.19) \quad \liminf_{t \rightarrow \infty} \int c d\nu_t \geq \rho^*.$$

Since this is true for any limit point ν of $\{\nu_t\}$ in $\mathcal{P}(\bar{\mathbf{S}} \times \mathbf{A})$ and for all sample points outside a set of probability zero, the desired result follows in this case also. \square

Remark 5.10. Some open problems arising in this context are:

(i) Can Assumption 5.20 be replaced by Assumption 5.19 while retaining the desired optimality?

(ii) If $\Pi_{SR} = \Pi_{SSR}$, will Assumption 5.19 hold automatically?

Remark 5.11. The condition in (5.19) implies a much stronger optimality, which will be discussed in §6.

Now, after the existence result of Theorem 5.11, an alternative treatment of the ACOE is possible. We will present a brief description without proofs. For details, see [24], [26], [28]. Define $h : \mathbf{S} \rightarrow \mathbb{R}$ by

$$(5.20) \quad h(i) = E_i^{f_0} \left[\sum_{t=0}^{\tau-1} (c(X_t, f_0(X_t)) - \rho^*) \right], \quad i \in \mathbf{S},$$

where $\tau = \min\{t \geq 1 : X_t = 0\}$ and $f_0 \in \Pi_{SD}$ is any sample path average cost optimal policy. In [22], [24], it is shown that $(h(\cdot), \rho^*)$ satisfies the ACOE under the following additional hypothesis called stability under local perturbations.

Assumption 5.21. Given an $f \in \Pi_{SSD}$ with $\rho(f) < \infty$, any $f' \in \Pi_{SD}$ obtained from f by changing the actions at most finitely many states is also stable and $\rho(f') < \infty$.

A sufficient, though not necessary, condition for Assumption 5.21 to hold is that every state has at most finitely many neighbors; i.e., for each $i \in \mathbf{S}$, there is a finite set $R_i \subset \mathbf{S}$ such that $P(j | i, \cdot) = 0$ for $j \notin R_i$.

In many cases, the solution (ρ^*, h) of the ACOE can be characterized (Theorem 5.12, below). The usual characterization of AC-optimal $f \in \Pi_{SD}$ in terms of the ACOE can also be proved for the foregoing.

THEOREM 5.12. *Assume Assumption 5.18 and let f_0, h be defined as above (cf. (5.20)). Let*

$$H = \{(\rho, w) : (\rho, w) \text{ satisfies the ACOE, } w(0) = 0, \inf w(\cdot) > -\infty\}.$$

Then (ρ^, h) is the unique element of H corresponding to the minimum value of ρ (i.e., if (ρ', w') is another element of H , then $\rho' \geq \rho^*$ with equality if and only if $w' = h$). Now, instead of Assumption 5.18, suppose that c is bounded and the following Lyapunov condition holds: There exists an $w : \mathbf{S} \rightarrow \mathbb{R}_+$, a finite $A \subset \mathbf{S}$ and an $\varepsilon > 0$ such that*

(a) $0 \in A$ and the set $\{i \in A^c : P(j | i, a) > 0, \text{ for some } j \in A, a \in \mathbf{A}\}$ is finite;

(b) $\lim_{i \rightarrow \infty} w(i) = \infty$;

(c) Under any $\pi \in \Pi, \mu \in \mathcal{P}(\mathbf{S})$

$$E_\mu^\pi [(w(X_{t+1}) - w(X_t) + \varepsilon)I\{X_t \notin A\} | \mathfrak{F}_t] \leq 0, \quad a.s.;$$

(d) There exists a random variable Z and a scalar $\lambda > 0$ such that $E[\exp(\lambda Z)] < \infty$ and, for all $b \geq 0$,

$$\mathcal{P}_\mu^\pi (|w(X_{t+1}) - w(X_t)| > b | \mathfrak{F}_t) \leq P(Z > b).$$

Then (ρ^, h) is the unique solution of the ACOE in the class $\{(\rho, w) : w(0) = 0, \limsup_{i \rightarrow \infty} h(i)/w(i) < \infty\}$.*

Remark 5.12. An alternative “intrinsic” formulation of the ACOE is also possible. For any $f \in \Pi_{SSD}$, define $h_f : \mathbf{S} \rightarrow \mathbb{R}$ by

$$h_f(i) = E_i^f \left[\sum_{t=0}^{\tau-1} (c(X_t, f(X_t)) - \rho(f)) \right], \quad i \in \mathbf{S}.$$

We say that f is *locally AC-optimal* if it yields a lower cost than any other element of Π_{SD} obtainable from f by changing f in at most finitely many states. In addition to the foregoing hypotheses, assume that every locally AC-optimal f is AC-optimal (for bounded c , a sufficient condition for this is that $\Pi_{SD} = \Pi_{SSD}$ and $\{\eta(f) : f \in \Pi_{SSD}\}$ is tight). We then have that f is sample path average cost optimal if and only if, for $i \in \mathbf{S}$,

$$h_f(i) = \inf_a \left\{ \sum_j P(j | i, a) h_f(j) + c(i, a) - \rho(f) \right\}.$$

This statement is “intrinsic” in the sense that all quantities (i.e., $h_f, \rho(f)$) are computable in terms of f . An interesting open problem is to characterize the most general conditions under which local AC-optimality implies AC-optimality.

Remark 5.13. The Lyapunov condition in Theorem 5.12(ii) implies Assumption 5.20 and has many other implications [26], but condition (ii)(d) there is rather strong, and, due to this, it may be difficult to construct such a function in a given situation. A partial answer to this question is given in [74]. It would be interesting to investigate if the Lyapunov conditions studied by [55], [91] (cf. Assumption 5.13), which do not involve condition (ii)(d) above, imply Assumption 5.20.

6. Borel state and action spaces. We consider in this section the case in which \mathbf{S} and \mathbf{A} are general Borel spaces. This is a natural setting for many problems, e.g., control of stock in water reservoirs, allocation of a resource between production and consumption, control of biological populations, harvesting a natural resource; see [17], [51], [82], and references therein for several examples. Also, the equivalent formulation of POCMP in terms of the conditional distribution of the (unobservable) state leads to a problem with an uncountable Borel state space, as we see in §7.

In this more general context, the ACOE is written as

$$(6.1) \quad \begin{aligned} \rho(x) + h(x) &= \inf_{a \in U(x)} \left\{ c(x, a) + \int_{\mathbf{S}} h(y) P(dy | x, a) \right\} \\ &= T(h)(x), \quad x \in \mathbf{S}, \end{aligned}$$

where $\rho, h \in \mathcal{M}(\mathbf{S})$. As in §5, a pair of functions (ρ, h) as above is called a solution to the ACOE, and, if ρ and h are bounded, we will say that the solution is bounded. Also, as in Theorem 5.1, our aim is to relate the AC problem to the existence of solutions to the ACOE. We have the following theorem.

THEOREM 6.1. *Suppose that (ρ, h) is a solution to the ACOE and that, for each policy $\pi \in \Pi_M$, the following holds:*

$$(6.2) \quad \lim_{t \rightarrow \infty} E_x^\pi \left[\frac{h(X_t)}{t} \right] = 0 \quad \forall x \in \mathbf{S}.$$

Then we have the following:

(i) *There holds*

$$(6.3) \quad \limsup_{n \rightarrow \infty} \frac{1}{n+1} E_x^\pi \left[\sum_{t=0}^n \rho(X_t) \right] \leq J(x, \pi),$$

and if $\pi \in \Pi_{SD}$ is such that $\pi(x)$ attains the infimum in (6.1), then equality is attained in (6.3);

(ii) *If $\rho(x) = \rho^* \in \mathbb{R}$, for all $x \in \mathbf{S}$, then $J^*(x) = \rho^*$, for all $x \in \mathbf{S}$, and any $\pi^* \in \Pi_{SD}$ such that $\pi^*(x)$ attains the infimum in (6.1) is average optimal.*

The proof of Theorem 6.1 follows that of Theorem 5.1 and is essentially contained in [177], more explicitly in [80], [148]; see also [78, pp. 66–68], [82, pp. 53–55], [150, pp. 93–94]. Note that (i), above, says that if $\rho(\cdot)$ is taken as the cost function to define another CMP $(\mathbf{S}, \mathbf{A}, U, P, \rho)$ then, for any $\pi \in \Pi_M$, the average cost incurred under the cost function $\rho(\cdot)$ does not exceed that under cost function $c(\cdot, \cdot)$.

Given the results above, it is of interest to find conditions under which there exists a solution (ρ, h) to the ACOE, satisfying (6.2). If h is bounded, then (6.2) is satisfied trivially. Also, if the random variables $\{h(X_t)\}$ are uniformly integrable under \mathcal{P}_x^π , for $\pi \in \Pi_M$ and $x \in \mathbf{S}$, then there exists a constant $0 < K_x^\pi < \infty$ such that $E_x^\pi [|h(X_t)|] \leq K_x^\pi$. Hence, if such a uniform integrability condition holds under *every* $\pi \in \Pi_M$ and $x \in \mathbf{S}$, then (6.2) is also satisfied trivially. The latter approach has been used by Shwartz and Makowski for some queueing problems [166]–[168].

6.1. Bounded costs. We first assume that $c(\cdot, \cdot)$ is bounded. When there are bounded solutions (ρ, h) to the ACOE, then much stronger results than those in Theorem 6.1 (i) can be obtained. To state these, some definitions are needed.

Let R and H be bounded, measurable, real-valued functions on \mathbf{S} , i.e., $R, H \in \mathcal{M}_b(\mathbf{S})$ and let $\pi^* \in \Pi$. Following the terminology of Dynkin and Yushkevich [51], the triplet (R, H, π^*) is said to be *canonical* if

$$(6.4) \quad J_N(x, \pi^*, H) = J_N^*(x, H) = H(x) + NR(x) \quad \forall N \in \mathbb{N}_0, \quad x \in \mathbf{S},$$

and $\pi^* \in \Pi$ is said to be a *canonical policy* if it is an element of some canonical triplet. Note that, if (R, H, π^*) is a canonical triplet, then π^* is N -stage optimal, for all $N \in \mathbb{N}_0$, when H is taken as the terminal cost. This concept was introduced by Yushkevich [204]. For finite models, Denardo and Fox [37] used a similar approach.

A policy $\pi^* \in \Pi$ is said to be *strong average optimal* if

$$(6.5) \quad \limsup_{N \rightarrow \infty} \frac{1}{N} J_N(x, \pi^*) \leq \liminf_{N \rightarrow \infty} \frac{1}{N} J_N(x, \pi) \quad \forall x \in \mathbf{S}, \pi \in \Pi.$$

Alternate definitions of strong average optimality are given in [69], [70]. Clearly, a strong average optimal policy π^* is also average optimal, and the limit of the sequence $\{1/N J_N(x, \pi^*)\}$, as $N \rightarrow \infty$, exists. An interpretation of (6.5) is that the “most pessimistic” average performance under π^* is no worse than the most “optimistic” performance under any other policy. We have the following result.

THEOREM 6.2. *Let $\pi^* \in \Pi_{SD}$, let $\rho, h \in \mathcal{M}_b(\mathbf{S})$, and let $c \in \mathcal{M}_b(\mathbf{K})$. Then (ρ, h, π^*) is a canonical triplet if and only if*

$$(6.6) \quad \rho(x) = \inf_{a \in U(x)} \left\{ \int_{\mathbf{S}} \rho(y) P(dy \mid x, a) \right\}$$

and

$$(6.7) \quad \rho(x) + h(x) = \inf_{a \in U(x)} \left\{ c(x, a) + \int_{\mathbf{S}} h(y) P(dy \mid x, a) \right\}$$

and $\pi^*(x)$ attains the infimum in both (6.6) and (6.7), for all $x \in \mathbf{S}$.

Proof. Necessity. Let (ρ, h, π^*) be a canonical triplet. Then, by (6.4),

$$(6.8) \quad \begin{aligned} h(x) + \rho(x) + N\rho(x) &= J_{N+1}^*(x, h) \\ &= T(J_N^*)(x) \\ &= c(x, \pi^*(x)) + \int_{\mathbf{S}} J_N^*(y, h) P(dy \mid x, \pi^*(x)). \end{aligned}$$

Since $J_0(x, \pi^*, h) = J_0^*(x, h) = h(x)$, then (6.7) follows from (6.8) by letting $N = 0$. Furthermore, since $\rho(\cdot)$, $h(\cdot)$, and $c(\cdot, \cdot)$ are bounded, then dividing both sides of (6.8) by N and letting $N \rightarrow \infty$ yields (6.6).

Sufficiency. Let (ρ, h) satisfy (6.6) and (6.7) and let $\pi^*(x)$ attain the infimum in these expressions. We use induction to show that (ρ, h, π^*) is a canonical triplet. For $N = 0$, this is trivially satisfied. Suppose that $N \in \mathbb{N}_0$ is the first integer for which (6.4) fails; then

$$\begin{aligned} J_N^*(x, h) &= T(J_{N-1}^*)(x) \\ &= T(h + (N-1)\rho)(x) \\ &= \inf_{a \in U(x)} \left\{ c(x, a) + \int_{\mathbf{S}} h(y) P(dy \mid x, a) + (N-1) \int_{\mathbf{S}} \rho(y) P(dy \mid x, a) \right\} \\ &\geq T(h)(x) + (N-1) \inf_{a \in U(x)} \left\{ \int_{\mathbf{S}} \rho(y) P(dy \mid x, a) \right\} \\ &= T(h)(x) + (N-1)\rho(x) = h(x) + N\rho(x). \end{aligned}$$

On the other hand,

$$\begin{aligned} J_N^*(x, h) &\leq J_N(x, \pi^*, h) \\ &= c(x, \pi^*(x)) + \int_{\mathbf{S}} J_{N-1}^*(y, \pi^*, h) P(dy \mid x, \pi^*(x)) \\ &= c(x, \pi^*(x)) + \int_{\mathbf{S}} [h(y) + (N-1)\rho(y)] P(dy \mid x, \pi^*(x)) \\ &= T(h)(x) + (N-1)\rho(x) = h(x) + N\rho(x) \end{aligned}$$

contradicting our hypothesis. Therefore, (ρ, h, π^*) is a canonical triplet. \square

The results in Theorem 6.2 were obtained by Yushkevich [204]; see also [51]. Note that (6.7) is the ACOE and that (6.6) allows $\rho(\cdot)$ to be treated as a constant, with respect to the optimization problem. Of course, if $\rho(x) = \rho^*$ for all $x \in \mathbf{S}$, then (6.6) is satisfied trivially. The coupled equations (6.6) and (6.7) were apparently introduced by Howard [95, pp. 61–62], in the context of finite state CMP for which, under some policies, $\{X_t\}$ has several ergodic classes, i.e., the so-called multichain case. In this case, different ergodic classes may have different optimal average cost, and $\rho(\cdot)$ gives this cost, as will be shown.

From Theorem 6.2, we see that the canonical policy π^* is a measurable selector for both (6.6) and (6.7). However, Assumption 2.2 in §2 is not enough to guarantee the existence of selectors in either (6.6) or (6.7), since ρ and h are assumed to be bounded and measurable functions, but not necessarily lower semicontinuous. For this situation, the following condition is needed.

Assumption 6.1. The transition kernel $P(\cdot | x, a)$ is *strongly continuous* in (x, a) ; that is, $u \in \mathcal{M}_b(\mathbf{S})$ implies that $\int_{\mathbf{S}} u(y)P(dy | \cdot, \cdot) \in C_b(\mathbf{K})$.

It follows that under Assumptions 2.1, 2.3, 6.1, measurable selectors exist for each of (6.6) and (6.7), and $\pi^* \in \Pi_{SD}$ will be a canonical policy if and only if it is a selector for both (6.6) and (6.7). If (ρ, h, π^*) is a canonical triplet, then (ρ, h) solves the ACOE, and (6.2) is satisfied, since h is bounded. Consequently, the results of Theorem 6.1 follow. The next result presents other important implications.

THEOREM 6.3. *Let (ρ, h, π^*) be a canonical triplet, and let $c \in \mathcal{M}_b(\mathbf{K})$. Then, for each $x \in \mathbf{S}$,*

- (i) $J_N(x, \pi^*) \leq J_N(x, \pi) + \text{span}(h)$, for every $\pi \in \Pi$;
- (ii) π^* is strong average optimal;
- (iii) $J(x, \pi^*) = J^*(x) = \rho(x)$;
- (iv) $h^-(x) + \rho(x)/(1 - \beta) \leq J_{\beta}^*(x) \leq h^+(x) + \rho(x)/(1 - \beta)$;
- (v) If $\rho(x) = \rho^* \in \mathbb{R}$, for all $x \in \mathbf{S}$, then, for every $\pi \in \Pi$,

$$\limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} c(X_t, A_t) \geq \rho^*, \quad \mathcal{P}_x^{\pi} \text{-a.s.},$$

when $X_0 = x$, and $\{A_t\}$ is generated using the policy π . Furthermore,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} c(X_t, A_t) = \rho^*, \quad \mathcal{P}_x^{\pi} \text{-a.s.},$$

if and only if

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} \Phi(X_t, A_t) = 0, \quad \mathcal{P}_x^{\pi} \text{-a.s.},$$

where $\Phi : \mathbf{K} \rightarrow \mathbb{R}$ is given by

$$\Phi(x, a) := c(x, a) + \int h(y)P(dy | x, a) - \rho^* - h(x);$$

- (vi) π^* is sample path average cost optimal.

Proof. To prove (i), note that, for all $\pi \in \Pi$,

$$\begin{aligned} J_N(x, \pi^*, h) &= E_x^{\pi^*} \left[\sum_{t=0}^{N-1} c(X_t, A_t) + h(X_N) \right] \\ &\leq E_x^{\pi} \left[\sum_{t=0}^{N-1} c(X_t, A_t) + h(X_N) \right] = J_N(x, \pi, h). \end{aligned}$$

Hence,

$$\begin{aligned} J_N(x, \pi^*) &\leq J_N(x, \pi) + E_x^{\pi} [h(X_N)] - E_x^{\pi^*} [h(X_N)] \\ &\leq J_N(x, \pi) + \text{span}(h) \quad \forall \pi \in \Pi. \end{aligned}$$

By the boundedness of $h(\cdot)$, we have that

$$\lim_{N \rightarrow \infty} \frac{1}{N} J_N(x, \pi^*, h) = \lim_{N \rightarrow \infty} \left[\frac{h(x) + N\rho(x)}{N} \right] = \rho(x).$$

Furthermore, since $J_N(x, \pi^*, h) = J_N(x, \pi^*) + E_x^{\pi^*}[h(X_N)]$, then

$$\rho(x) = \lim_{N \rightarrow \infty} \frac{1}{N} J_N(x, \pi^*),$$

and (ii)–(iii) follows from (i).

Next, since (ρ, h) solve the ACOE, then (ρ, h^-) and (ρ, h^+) are also solutions to the ACOE. Since $h^-(\cdot) \leq 0 \leq h^+(\cdot)$, then by Lemma 2.1 we have that $T(h^-) \leq T(\beta h^-) = T_\beta(h^-)$, and $T(h^+) \geq T(\beta h^+) = T_\beta(h^+)$. Then, (iv) follows by induction, using Theorem 2.1 (iv); see [64].

Turning our attention to (v) and (vi), observe that, due to (6.7), $\Phi(x, a) \geq 0$ for all $(x, a) \in \mathbf{K}$. Also, by the (Markov) property (2.3) in §2, we have that, for any $\pi \in \Pi$,

$$\Phi(X_t, A_t) = E_x^\pi \left[c(X_t, A_t) + h(X_{t+1}) - \rho^* - h(X_t) \mid H_t, A_t \right], \quad \mathcal{P}_x^\pi\text{-a.s.}$$

Let

$$Z_t := c(X_t, A_t) + h(X_{t+1}) - h(X_t) - \rho^* - \Phi(X_t, A_t)$$

and

$$M_N := \sum_{t=0}^{N-1} Z_t = \sum_{t=0}^{N-1} c(X_t, A_t) - N\rho^* + h(X_N) - h(X_0) - \sum_{t=0}^{N-1} \Phi(X_t, A_t).$$

Note that $\{Z_t\}$ is a $(\mathfrak{G}_t, \mathcal{P}_x^\pi)$ martingale difference, where $\mathfrak{G}_t := \sigma(H_{t+1}, A_{t+1})$. Since $\{Z_t\}$ is bounded uniformly in t , by the martingale stability theorem

$$\lim_{N \rightarrow \infty} \frac{M_N}{N} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} Z_t = 0, \quad \mathcal{P}_x^\pi\text{-a.s.}$$

Therefore, by the boundedness of $h(\cdot)$,

$$\lim_{N \rightarrow \infty} \left[\frac{1}{N} \sum_{t=0}^{N-1} c(X_t, A_t) - \rho^* - \frac{1}{N} \sum_{t=0}^{N-1} \Phi(X_t, A_t) \right] = 0, \quad \mathcal{P}_x^\pi\text{-a.s.}$$

Finally, (v) and (vi) follow, since $\Phi(x, a) \geq 0$ for all $(x, a) \in \mathbf{K}$ and since, for a canonical policy π^* , $\Phi(X_t, A_t) = 0$, $\mathcal{P}_x^{\pi^*}$ -a.s. \square

The results in (i)–(iii) of Theorem 6.3 are essentially contained in [51, Chap. 7]; that in (iv) is motivated by similar results in [136] and [64]; (v) and (vi) are due to Georgin [72], see also [82, pp. 52–55]. Also, the function Φ defined in (v) was introduced by Mandl [124] and is often referred to as *Mandl's discrepancy function*.

In view of Theorem 6.3, it follows that a canonical triplet yields the desired results. We therefore look for conditions on the primary objects like the cost function c and transition kernel P , which imply the existence of a canonical triplet, so that the theory

can be used in a given practical situation. To this end, a standard procedure is to assume some ergodicity conditions that will ensure the existence of a bounded solution to the ACOE. We have already discussed several such conditions for the countable state case (cf. Assumptions 5.1–5.5). Analogues of such assumptions are also available in the literature, an extensive survey of which appears in [86]. We will focus on a particular ergodicity condition that not only subsumes many such conditions but also facilitates easily implementable numerical schemes. Our presentation here follows essentially that in [82, Chap. 3].

Assumption 6.2. There exists a number $\alpha < 1$ such that

$$\sup_{k,k' \in \mathbf{K}} \|P(\cdot | k) - P(\cdot | k')\|_{TV} \leq 2\alpha,$$

where $\|\cdot\|_{TV}$ denotes the total variation norm.

Example 6.1. Let $\mathbf{S} = \mathbb{R}$, $\mathbf{A} \subset \mathbb{R}$, a compact set. Consider the system

$$X_{t+1} = F(X_t, A_t) + G(X_t)W_t, \quad X_0 = x,$$

where $F : \mathbb{R} \times \mathbf{A} \rightarrow \mathbb{R}$, $G : \mathbb{R} \rightarrow \mathbb{R}$ are bounded, continuous and $G(\cdot) > 0$, and $\{W_t\}$ is a sequence of independent $N(0, 1)$ random variables ($N(a, b)$ stands for the Gaussian distribution with mean a and variance b). In this case, the transition kernel is given by

$$P(\cdot | x, a) = N(F(x, a), G^2(x)).$$

Using the assumed conditions on F, G we can show that Assumption 6.2 holds. We omit the details. An important consequence of Assumption 6.2 is given below; for a proof and further discussion, see [82, Chap. 3].

LEMMA 6.1. *Suppose that Assumption 6.2 holds. Then, for any $f \in \Pi_{SD}$, the corresponding process $\{X_t\}$ has a unique invariant measure $\eta(f) \in \mathcal{P}(\mathbf{S})$ satisfying*

$$(6.9) \quad \|P^t(\cdot | x, f(x)) - \eta(f)(\cdot)\|_{TV} \leq 2\alpha^t, \quad t = 0, 1, \dots,$$

where $P^t(\cdot | x, f(x))$ denotes the t -step transition probability measure under f with $X_0 = x$.

Remark 6.1. (a) Lemma 6.1 also holds for any $f \in \Pi_{SR}$.

(b) It follows from (6.9) that, for any $f \in \Pi_{SD}$, $P^t(\cdot | x, f(x))$ converges to $\eta(f)$ in total variation norm, uniformly in x , and at a geometric rate.

(c) It is clear that, for any $f \in \Pi_{SD}$,

$$J(\mu, f) = \int_{\mathbf{S}} c(x, f(x))\eta(f)(dx)$$

for any initial law μ .

(d) Compare (6.9) with Assumption 5.5 In view of Theorem 5.4, Assumption 6.2 may be viewed as a representative counterpart of Assumptions 5.1–5.5 for the general state space case.

We now introduce the concept of span-contraction.

DEFINITION 6.1. Let $T : \mathcal{M}_b(\mathbf{S}) \rightarrow \mathcal{M}_b(\mathbf{S})$. T is said to be a *span-contraction* if, for some $\gamma \in [0, 1)$,

$$\text{span}(Tu - Tv) \leq \gamma \text{span}(u - v) \quad \text{for all } u, v \in \mathcal{M}_b(\mathbf{S}).$$

Let \sim be the equivalence relation on $\mathcal{M}_b(\mathbf{S})$ defined by $u \sim v$ if and only if there exists some constant C such that $u(x) - v(x) = C$ for all $x \in \mathbf{S}$. Let $\widetilde{\mathcal{M}}_b(\mathbf{S}) = \mathcal{M}_b(\mathbf{S}) / \sim$, the quotient space, endowed with the quotient norm induced by the span seminorm. For $v \in \mathcal{M}_b(\mathbf{S})$, let \tilde{v} denote the corresponding element of $\widetilde{\mathcal{M}}_b(\mathbf{S})$ and $\widetilde{T} : \widetilde{\mathcal{M}}_b(\mathbf{S}) \rightarrow \widetilde{\mathcal{M}}_b(\mathbf{S})$ be the canonically induced map, i.e., $\widetilde{T}\tilde{v} = \widetilde{Tv}$, $v \in \mathcal{M}_b(\mathbf{S})$. It is easily seen that, if T is a span-contraction on $\mathcal{M}_b(\mathbf{S})$, then \widetilde{T} is a contraction on $\widetilde{\mathcal{M}}_b(\mathbf{S})$ and therefore has a unique fixed point. In turn, it follows that the map T has a span-fixed point; i.e., there exists a $v^* \in \mathcal{M}_b(\mathbf{S})$ such that $\text{span}(Tv^* - v^*) = 0$ or, equivalently, $Tv^* - v^*$ is a constant. It also follows that any two span-fixed points of T must differ by a constant.

We now replace Assumption 2.3 with the following

Assumption 6.3. (i) The multifunction $U(x)$ is continuous; (ii) $c(\cdot, \cdot) \in C_b(\mathbf{K})$.

We have the following result; for a proof, see [82, Lemma 3.5].

LEMMA 6.2. *Under Assumptions 2.2, 6.2, and 6.3, the operator T defined in (2.5) maps $C_b(\mathbf{S})$ to $C_b(\mathbf{S})$ and is a span-contraction.*

COROLLARY 6.1. *Under Assumptions 2.2, 6.2, and 6.3, the ACOE has a bounded solution $(\rho^*, h^*) \in \mathbb{R} \times C_b(\mathbf{S})$.*

Proof. This follows from the fact that there exists a $h^* \in C_b(\mathbf{S})$ such that $\text{span}(Th^* - h^*) = 0$. Hence, $Th^* = h^* + \rho^*$ for some constant ρ^* . \square

Remark 6.2. (a) Assume Assumptions 6.2 and 6.3. Let $(\rho^*, h^*) \in \mathbb{R} \times C_b(\mathbf{S})$ be a solution to the ACOE and fix $x_0 \in \mathbf{S}$. Define $h(\cdot) = h^*(\cdot) - h^*(x_0)$. Then (ρ^*, h) is also a solution to the ACOE. By the span-contraction property of T , it is the unique solution in $\mathbb{R} \times C_b(\mathbf{S})$ satisfying $h(x_0) = 0$; i.e., if $(\rho', h') \in \mathbb{R} \times C_b(\mathbf{S})$ is any other solution of the ACOE in $\mathbb{R} \times C_b(\mathbf{S})$ such that $h'(x_0) = 0$, then $\rho' = \rho$ and $h' = h$.

(b) In view of the span-contraction property of the operator T , the value iteration scheme described in §4 can be extended to this case; for details, we refer to [82, Chap. 3].

(c) Note that Corollary 6.1 asserts the existence of a canonical triplet.

Remark 6.3. In §4 we have identified the duality between the linear programming formulation and the ACOE under the irreducibility assumption. This has been extended by Yamada [203] to the case when the state space \mathbf{S} is a compact subset of \mathbb{R}^n and the transition law has a density that satisfies a certain ‘‘positivity’’ condition. Hernández-Lerma, Hennet, and Lasserre [84] have further extended this result to the Borel state space setting under Assumption 6.2.

Kurano [105]–[107] has studied the problem for compact state and action spaces, under the hypothesis of Doeblin. Doeblin’s condition for the general state space can be described as follows.

Assumption 6.4. There exists a nontrivial finite measure μ on $(\mathbf{S}, \mathcal{B}(\mathbf{S}))$, a positive integer ℓ , and an $\varepsilon > 0$ such that

$$P^\ell(A \mid x, f(x)) \geq 1 - \varepsilon \quad \text{if } \mu(A) \geq \varepsilon,$$

for all $f \in \Pi_{SD}$ and $x \in \mathbf{S}$.

THEOREM 6.4. *Let the state and action spaces be compact and Assumptions 6.3 and 6.4 hold. Then there exist an $f \in \Pi_{SD}$ and a set $A \in \mathcal{B}(\mathbf{S})$ with $\mu(A) > \varepsilon$ such that $P(A \mid x, f(x)) = 1$ for all $x \in \mathbf{S}$, and f is optimal, provided that the initial law is supported on the set A .*

Furthermore, assume the following.

Assumption 6.5 (reachability). For any $x \in S$ and $D \in \mathcal{B}(S)$ with $\mu(D) > \varepsilon$ (μ and ε as in Assumption 6.4, there exists a $\pi \in \Pi$ such that

$$P_x^\pi \left(\bigcup_{t=0}^\infty \{X_t \in D\} \right) = 1.$$

Assumption 6.6. One of the following two conditions is satisfied:

- (i) $\mu(\partial D) = 0$ if $\mu(D) > 0$, where ∂D denotes the boundary of D ;
- (ii) For each $D \in \mathcal{B}(S)$ with $\mu(D) > \varepsilon$, $P(D \mid x, a)$ is continuous in (x, a) .

THEOREM 6.5. *Under Assumptions 6.3–6.6 there exists an $f \in \Pi_{SD}$, which is optimal.*

Remark 6.4. (a) The proof of Theorem 6.3 exploits the idea involved in Lemma 5.1 of extracting a stationary randomized policy from a limit point of empirical processes. A novel idea in [105] is to remove the randomization by using the ergodic decomposition of Markov processes under Assumption 6.4. The compactness is used to ensure the tightness of the empirical processes under any policy. This can be dropped if the cost function has a penalizing condition or if there is a blanket stability of Lyapunov type. The details closely mimic the development at the end of §5.

(b) Wijngaard [201] has also obtained the existence of an optimal $f \in \Pi_{SD}$ under Doeblin’s condition using an operator theoretic method.

We will now discuss the vanishing discount approach to obtain a bounded solution to the ACOE. For a fixed $x_0 \in S$, let $h_\beta(\cdot) = J_\beta^*(\cdot) - J_\beta^*(x_0)$ denote the differential discounted value function. For a general state space, the usual diagonalization procedure used on a countable state space is not amenable. Nevertheless, if $h_\beta(\cdot)$ is uniformly bounded and equicontinuous, then we can use a more subtle diagonalization involving the Arzela–Ascoli theorem to take the required limits and obtain a bounded solution to the ACOE. This was studied by Ross [148]. Following [17], [72], [73], we will discuss some sufficient conditions to obtain the required uniform boundedness and equicontinuity of $h_\beta(\cdot)$.

Assumption 6.7. For each $\beta \in (\beta', 1)$, for some $0 < \beta' < 1$, and $f_\beta \in \Pi_{SD}$, the corresponding state process has a unique invariant probability measure $\eta(f_\beta)$ such that

$$(6.10) \quad \sup_{\substack{x \in S \\ \beta \in (\beta', 1)}} \sum_{t=1}^\infty \|P^t(\cdot \mid x, f_\beta(x)) - \eta(f_\beta)(\cdot)\|_{TV} < \infty.$$

The following result is now easy to establish.

LEMMA 6.3. *Under Assumptions 6.1, 6.3, and 6.7, $h_\beta(\cdot) := J_\beta^*(\cdot) - J_\beta^*(x_0)$, $x_0 \in S$ fixed, is uniformly bounded, and is equicontinuous for $\beta \in (\beta', 1)$.*

COROLLARY 6.2. *Under Assumptions 6.1, 6.3, and 6.7, the ACOE has a solution (ρ^*, h) such that $h \in C_b(S)$.*

Remark 6.5. If Assumption 6.4 is satisfied and we further impose the condition that, for every $f \in \Pi_{SD}$, the corresponding state process has a single ergodic class, then (6.10) holds. In particular, if $P(dy \mid x, a)$ has a density $p(y, x, a)$, with respect to some σ -finite measure μ , and there exists a nonnegative measurable function p_0 satisfying $\int p_0(y)\mu(dy) > 0$ and $p(y, x, a) \geq p_0(y)$, for all (x, a) , then Assumption 6.4 holds and (6.10) can be easily verified. If $(x, a) \rightarrow p(y, x, a)$ is continuous, then by Scheffe’s theorem, $p(\cdot \mid x, a)$ is strongly continuous in (x, a) .

6.2. Unbounded costs. We now drop the boundedness condition on the cost function and discuss some recent developments involving refinements and extensions of the vanishing discount approach. Since for unbounded costs the uniform boundedness of the differential discounted value function $h_\beta(\cdot)$ is rather unnatural, we attempt to extend the procedure of [155], [156] to the present case. To this end, we make the following analogues of Assumptions 5.14–5.16.

Assumption 6.8. There exists a nonnegative function $b \in \mathcal{M}(\mathbf{S})$, a constant $M \geq 0$, and a sequence $\{\beta_n\} \subset (0, 1)$, $\beta_n \uparrow 1$, such that for all $x \in \mathbf{S}$, (i) $-M \leq h_{\beta_n}(x) \leq b(x)$, and (ii) $\int_{\mathbf{S}} b(y)P(dy | x, a) < \infty$, for all $a \in U(x)$.

Assumption 6.9. There exists a policy π and an initial state \hat{x} such that $J(\hat{x}, \pi) < \infty$.

Assumption 6.10. There exists $\beta' \in (0, 1)$ such that $\sup_{\beta \in (\beta', 1)} \tilde{h}_\beta(x) < \infty$, where $\tilde{h}_\beta(x) = J_\beta^*(x) - \inf_{x \in \mathbf{S}} J_\beta^*(x)$.

Assumption 6.11. The transition kernel $P(\cdot | x, a)$ is strongly continuous in a , for each $x \in \mathbf{S}$.

Under Assumptions 6.8 and 6.11, defining $h(x) = \liminf_{n \rightarrow \infty} h_{\beta_n}(x)$, $x \in \mathbf{S}$, and using Fatou's lemma, we can show that there exists a constant ρ^* such that

$$(6.11) \quad \lim_{n' \rightarrow \infty} (1 - \beta_{n'})J_{\beta_{n'}}^*(x) = \rho^* \quad \text{for all } x \in \mathbf{S},$$

where $\beta_{n'} \uparrow 1$ is a subsequence of $\{\beta_n\}$, and

$$(6.12) \quad \rho^* + h(x) \geq \min_{a \in U(x)} \left\{ c(x, a) + \int_{\mathbf{S}} h(y)P(dy | x, a) \right\}, \quad x \in \mathbf{S},$$

which is the ACOI (see (5.16)) for this case. Similarly, under Assumptions 6.9–6.11, we can find a constant ρ^* such that, along a suitable sequence $\beta_n \in (\beta', 1)$, $\beta_n \uparrow 1$, $\lim_{n \rightarrow \infty} (1 - \beta_n) \inf_{x \in \mathbf{S}} J_\beta^*(x) = \rho^*$. Then, defining $h(x) = \liminf_{n \rightarrow \infty} \tilde{h}_{\beta_n}(x)$, we can deduce (6.12). Thus, we have the following result.

THEOREM 6.6. *Under Assumptions 6.8 and 6.11 or under Assumption 6.9–6.11, there exists a constant ρ^* and a function h , which is bounded below and satisfies (6.12). Any policy $\pi \in \Pi_{SD}$ realizing the minimum on the right-hand side of (6.12) is average optimal and ρ^* is the minimum average cost.*

Remark 6.6. For details, we refer to [83], [85], [140]. In the case of a countable state space, a number of sufficient conditions on the transition kernel and the cost function that enable us to verify Assumptions 5.14–5.16 are available, as mentioned in §5. This does not seem to be the case for a general Borel state space model, although several interesting examples have been studied in [83], [85], and [140]. Also, Assumption 6.11 is a very strong condition and will not, in general, be satisfied for the transition kernel of the equivalent problem for a partially observable model. Thus, this case needs further investigation. Finally, note that Assumption 6.10 may in principle be easier to verify than Assumption 6.8.

Remark 6.7. We note that Theorem 6.6 provides only an ACOI, and not the ACOE. In many situations, the discounted value function is convex (e.g., in linear systems with quadratic cost [14]), or concave (e.g., the separated problem in partially observable models). This class of problems has been used in [61] to obtain the ACOE under Assumptions 6.8, 6.11, and some additional assumptions.

7. Partially observable controlled Markov processes. Thus far, we have assumed that the complete history of the process H_t is available to the decision-maker, at each stage $t \in T$. However, in many situations, some components of the state process may not be directly available to the controller since, e.g., it may be impossible or too costly to measure these. Furthermore, due to imprecisions in the measuring devices, only noisy observations of the state may be available. When these situations arise, the problem is said to be a partially observable controlled Markov process. Here, we study POCMP with finite or countably infinite state and observation spaces, and finite or compact action set. A major portion of our exposition concentrates on the vanishing discount method, where we see that the particular structure of the POCMP can be employed to yield stronger results than those available for general Borel spaces. We also review Borkar's convex analytic approach, specialized to the partially observable case [26].

7.1. Models with partial state information. The model for this problem is essentially that in [51, Chap. 8] and is as follows. The state process is described by a pair $\{X_t, Y_t\}_{t \in T}$ taking values in a product of Borel spaces $\mathbf{X} \times \mathbf{Y}$. Only the second component $\{Y_t\}_{t \in T}$ of the state process is available for decision-making, and, reflecting this, \mathbf{Y} is called the *observation or message space*, and Y_t the *observation process*. With \mathbf{A} denoting the action space, the evolution of the system is governed by a measurable stochastic kernel P on $\mathbf{X} \times \mathbf{Y}$ given $\mathbf{X} \times \mathbf{Y} \times \mathbf{A}$.

Let $\mu \in \mathcal{P}(\mathbf{X} \times \mathbf{Y})$ be an initial distribution of the state. Decomposing (disintegrating) the measure μ , we have

$$\mu(dx, dy) = \bar{Q}_0(dy) \psi_0(dx | y),$$

where \bar{Q}_0 is the marginal of μ on \mathbf{Y} and ψ_0 is a version of the regular conditional law, defined \bar{Q}_0 almost surely; we pick any version from this equivalence class and keep it fixed thereafter. Note that knowledge of μ , since the value of Y_0 is available to the controller, implies that an a posteriori distribution ψ_0 (given $Y_0 = y$) for the unobserved initial state is introduced. We include ψ_0 into the *observed history* by letting

$$\bar{\mathbf{H}}_0 := \mathcal{P}(\mathbf{X}) \times \mathbf{Y}, \quad \bar{\mathbf{H}}_t := \bar{\mathbf{H}}_{t-1} \times \mathbf{Y} \times \mathbf{A}, \quad t \in \mathbb{N}.$$

The set of admissible actions is specified by a strict, measurable, compact-valued multifunction $U : \mathbf{Y} \rightarrow \mathcal{B}(\mathbf{A})$. Hence, in this context, an admissible policy is a sequence $\pi = \{\pi_t\}_{t \in T}$ of Borel measurable stochastic kernels π_t on \mathbf{A} given $\bar{\mathbf{H}}_t$ satisfying, for all $t \in T$, the constraint

$$\pi_t(U(y_t) | \bar{h}_t) = 1 \quad \forall \bar{h}_t \in \bar{\mathbf{H}}_t.$$

The set of all admissible policies is again denoted by Π .

Remark 7.1. In general, decisions take into account past and present information, not just the last observation. Note that the constraints on the actions cannot depend on the unobservable component X_t of the state. If this type of constraint must be included in the model, then it must be provided to the controller as an additional observation. Similarly, if the cost process $\{c(X_t, Y_t, A_t)\}$ is available to the controller, then it should also be regarded as an additional component in the observation process [51, p. 201].

Remark 7.2. Quite often, μ is specified as

$$\mu(dx, dy) = Q_0(dy | x) \mu_0(dx),$$

where $\mu_0 \in \mathcal{P}(\mathbf{X})$ is an initial distribution for X_0 , and Q_0 is a stochastic kernel on \mathbf{Y} given \mathbf{X} [15, Chap. 10], [82, Chap. 4].

With $\mu \in \mathcal{P}(\mathbf{X} \times \mathbf{Y})$ and an admissible policy π specified, there exists a unique probability measure \mathcal{P}_μ^π on $(\Omega, \mathcal{B}(\Omega))$, where $\Omega := (\mathbf{X} \times \mathbf{Y} \times \mathbf{A})^\infty$, defined by

$$\begin{aligned} \mathcal{P}_\mu^\pi(dx_0, dy_0, da_0, \dots, da_{t-1}, dx_t, dy_t) \\ = \mu(dx_0, dy_0) \pi_0(da_0 | \psi_0, y_0) P(dx_1, dy_1 | x_0, y_0, a_0) \cdots \\ \pi_{t-1}(da_{t-1} | \psi_0, y_0, a_0, \dots, y_{t-1}) P(dx_t, dy_t | x_{t-1}, y_{t-1}, a_{t-1}). \end{aligned}$$

7.2. Transformation into a completely observable model. A common approach in the analysis of a partially observable (PO) model is to construct a completely observable (CO) model, equivalent to the original one in the sense that corresponding policies have equal costs. The advantages in doing this are obvious, since the theory of CO problems is much better developed. However, the price usually paid is that the dimensionality of the new state space is substantially larger than that of the original one.

Such an equivalent CO problem can be obtained in many ways. The main idea is to specify an *information state process* that summarizes, at each time, all relevant information for decision-making. Clearly, $\bar{H}_t = (\psi_0, Y_0, A_0, \dots, A_{t-1}, Y_t)$ can be used as an information state process, but this leads to a nonstationary CO model, in which “growing memory” difficulties arise; see [15, Chap. 10]. We present here the more standard approach where the inferential knowledge of X_t is summarized using its conditional probability distribution, given the entire observed history up to time t . We first present the construction of the equivalent CO model for general Borel state spaces and then specialize to models with countable state space. Also, the following assumption will be in effect throughout this section.

Assumption 7.1. The transition kernel $P(\cdot | x, y, a)$ and the cost function $c(x, y, a)$ do not depend on y , and $U(y) = \mathbf{A}$ for all $y \in \mathbf{Y}$.

7.2.1. Borel state space. Given a PO model $(\mathbf{X} \times \mathbf{Y}, \mathbf{A}, U, P, c)$ satisfying Assumption 7.1, we construct a CO model $(\mathcal{P}(\mathbf{X}), \mathbf{A}, \tilde{U}, \mathcal{K}, \tilde{c})$ as follows. Let $\{\Psi_t, Y_t\}_{t \in T}$ and $\{\tilde{H}_t\}_{t \in T}$ denote the state process and the history spaces, respectively. The set of admissible actions is selected by letting $\tilde{U}(\psi) = \mathbf{A}$ for all $\psi \in \mathcal{P}(\mathbf{X})$. We define the cost function \tilde{c} by

$$(7.1) \quad \tilde{c}(\psi, a) := \int_{\mathbf{X}} c(x, a) \psi(dx), \quad \psi \in \mathcal{P}(\mathbf{X}).$$

It remains to construct the transition kernel \mathcal{K} . Working on the canonical sample space $\tilde{\Omega} = (\mathcal{P}(\mathbf{X}) \times \mathbf{A})^\infty$, we first define a stochastic kernel q on $\mathbf{X} \times \mathbf{Y}$ given $\mathcal{P}(\mathbf{X}) \times \mathbf{A}$ by

$$(7.2) \quad q(dx, dy | \psi, a) := \int_{\mathbf{X}} P(dx, dy | x', a) \psi(dx'), \quad \psi \in \mathcal{P}(\mathbf{X}),$$

and, decomposing q , we obtain

$$(7.3) \quad q(dx, dy | \psi, a) = Q(dy | \psi, a) \Psi(dx | \psi, a, y).$$

Equation (7.3) is the *filtering equation*. For fixed (ψ, a) , the map $y \mapsto \Psi$, as defined implicitly in (7.3), is a measurable mapping from \mathbf{Y} to $\mathcal{P}(\mathbf{X})$. Consequently, along

with the distribution Q on \mathbf{Y} , it induces a distribution \mathcal{K} on $\mathcal{B}(\mathcal{P}(\mathbf{X}))$, which is a measurable function of (ψ, a) or, in other words, a stochastic kernel on $\mathcal{P}(\mathbf{X})$ given $\mathcal{P}(\mathbf{X}) \times \mathbf{A}$. It follows that the model $(\mathcal{P}(\mathbf{X}), \mathbf{A}, \mathcal{K}, \tilde{c})$, with state process $\{\Psi_t\}_{t \in T}$, forms a completely observable controlled Markov process, with transition kernel given by

$$(7.4) \quad \mathcal{K}(B \mid \psi, a) := \int_{\mathbf{Y}} I\{\Psi(\cdot \mid \psi, a, y) \in B\} Q(dy \mid \psi, a), \quad B \in \mathcal{B}(\mathcal{P}(\mathbf{X})).$$

The distribution $\tilde{\mu}_0$ of Ψ_0 , corresponding to an initial distribution μ of the PO model, is taken to be

$$(7.5) \quad \tilde{\mu}_0(B) := \int_{\mathbf{Y}} \mu(B, dy), \quad B \in \mathcal{B}(\mathcal{P}(\mathbf{X})).$$

Given a history $\bar{h}_t = (\psi_0, y_0, \dots, a_{t-1}, y_t) \in \bar{\mathbf{H}}_t$ in the PO model, we can construct ψ_1, ψ_2, \dots in a recursive manner by starting from ψ_0 and, having obtained ψ_{t-1} , solving for Ψ in (7.3), with $(\psi, a, y) = (\psi_{t-1}, a_{t-1}, y_t)$, and letting $\psi_t = \Psi$. In this manner, we obtain a corresponding history $\tilde{h}_t = (\psi_0, a_0, \dots, a_{t-1}, \psi_t) \in \tilde{\mathbf{H}}_t$ for the CO model; we denote this correspondence by the map $g_t : \bar{\mathbf{H}}_t \rightarrow \tilde{\mathbf{H}}_t$. We can then assign to each admissible policy $\tilde{\pi} \in \tilde{\Pi}$ in the CO model a corresponding policy $\pi = g^*(\tilde{\pi})$ in the PO model, defined by

$$(7.6) \quad \pi_t(\cdot \mid \bar{h}_t) := \tilde{\pi}_t(\cdot \mid g_t(\bar{h}_t)), \quad \bar{h}_t \in \bar{\mathbf{H}}_t.$$

Clearly, every policy $\pi \in \Pi$ can also be regarded as a policy in $\tilde{\Pi}$; in other words, the map g^* is onto. If $\mathcal{P}_{\tilde{\mu}}^{\tilde{\pi}}$ is the probability measure induced by the policy $\tilde{\pi}$ and the initial distribution $\tilde{\mu}$ (corresponding to μ) on the canonical sample space $\tilde{\Omega}$, then, for each $C \in \mathcal{B}(\mathbf{X})$,

$$(7.7) \quad \mathcal{P}_{\tilde{\mu}}^{g^*(\tilde{\pi})}(X_t \in C \mid \bar{\mathbf{H}}_t = \bar{h}_t) = \Psi_t(C), \quad \mathcal{P}_{\tilde{\mu}}^{\tilde{\pi}}\text{-a.s.}$$

Utilizing (7.1), (7.4), and (7.5), it can be verified that

$$(7.8) \quad E_{\tilde{\mu}}^{g^*(\tilde{\pi})}[c(X_t, A_t)] = E_{\tilde{\mu}}^{\tilde{\pi}}[\tilde{c}(\Psi_t, A_t)] \quad \forall t \in T,$$

thus establishing that the two models are indeed equivalent as claimed. It follows that the process Ψ_t summarizes all information, relevant for control purposes, and is called for this purpose a *sufficient statistic* (see [50], [161], [162]). We define the set of *separated policies* Π_S as those policies $\pi \in \Pi$ for which there a Markov policy $\tilde{\pi}$ on the equivalent CO problem such that $\pi = g^*(\tilde{\pi})$, as defined in (7.6). In other words, with $\tilde{\pi} = \{f_t\}_{t \in T} \in \tilde{\Pi}_M$, $f_t : \mathcal{P}(\mathbf{X}) \rightarrow \mathcal{P}(\mathbf{A})$ and for each initial distribution $\mu \in \mathcal{P}(\mathbf{S})$,

$$\pi_t(\cdot \mid \bar{h}_t) = f_t(\Psi_t)(\cdot), \quad \mathcal{P}_{\tilde{\mu}_0}^{\tilde{\pi}}\text{-a.s.}$$

Thus, the actions taken using a separated policy only depend on $\bar{\mathbf{H}}_t$ through the conditional distribution of X_t . In other words, the following *separation principle* holds: If an optimal policy exists in Π , one exists in Π_S . Hence, the process can be controlled optimally by first estimating the state via the conditional distribution and

choosing control actions based solely on the latter. These and other results, in various degrees of generality, were independently obtained by various authors, e.g., [3], [5], [89], [138], [151], [163], [174], [175], [199], [205].

Example 7.1. A partially observable version of the stochastic nonlinear system in Example 2.1 is described by the equations

$$\begin{aligned} X_{t+1} &= F(X_t, A_t, W_t), \\ Y_t &= G(X_t, A_{t-1}, V_t), \\ Y_0 &= G_0(X_0, V_0), \end{aligned}$$

where G and G_0 are Borel measurable, and the disturbance $\{V_t\}_{t \in T}$ is an i.i.d. sequence of random variables taking values in a Borel space \mathbf{V} , with a common distribution \mathcal{P}_V ; furthermore, it is assumed that X_0 , $\{W_t\}$, and $\{V_t\}$ are mutually independent.

7.2.2. Countable state space. We now specialize to the case where the state space $\mathbf{X} \times \mathbf{Y}$ is a finite or countably infinite set, the action space \mathbf{A} is a finite or compact set and with Assumption 7.1 in effect. Thus, $U(y) = \mathbf{A}$ for all $y \in \mathbf{Y}$, and the kernel of the process takes the form $P(x', y' | x, a)$. We also assume that the cost c and the kernel P are continuous with respect to $a \in \mathbf{A}$. The space $\mathcal{P}(\mathbf{X})$ is identified with the set Δ of probability vectors, i.e.,

$$(7.9) \quad \Delta := \left\{ \psi \in [0, 1]^{\mathbf{X}} : \sum_{x \in \mathbf{X}} \psi(x) = 1 \right\}$$

endowed with the topology given by the metric

$$d(\psi_1, \psi_2) := \sum_{x \in \mathbf{X}} |\psi_1(x) - \psi_2(x)| = \|\psi_1 - \psi_2\|_1,$$

where $\|\cdot\|_1$ stands for the standard ℓ_1 -norm on $\mathbb{R}^{\mathbf{X}}$.

In general, the recursive (filtering) equation (7.3) used to compute ψ_{t+1} , is obtained via a decomposition of measures technique; see [15, Chap. 10], [51, Chap. 8], [82, Chap. 4], [205]. This is particularly simple to accomplish (using the Bayes rule) when \mathbf{X} and \mathbf{Y} are countable or when the system is described by a linear system function and the disturbances are Gaussian; see [5], [14], [103], [174], [175]. For this purpose, we need the following definitions (compare with (7.2), (7.3)):

$$(7.10) \quad q(x, y | \psi, a) := \sum_{x' \in \mathbf{X}} P(x, y | x', a) \psi(x'),$$

$$(7.11) \quad V(y, \psi, a) := \sum_{x \in \mathbf{X}} q(x, y | \psi, a),$$

$$(7.12) \quad T(y, \psi, a)(\cdot) := \begin{cases} \frac{q(\cdot, y | \psi, a)}{V(y, \psi, a)}, & \text{if } V(y, \psi, a) \neq 0, \\ 0, & \text{otherwise.} \end{cases}$$

Note that the map $\psi \rightarrow T(y, \psi, a)$ maps Δ into itself. In the countable case, ψ_t can be computed by letting $\psi_t = T(y_t, \psi_{t-1}, a_{t-1})$. Here, $V(y, \psi, a)$ is interpreted as the (one-step ahead) conditional probability of the observation being y given an a priori distribution ψ for the core state, under decision a . Likewise, $T(y, \psi, a)$ is interpreted

as the a posteriori conditional probability distribution of the core state, given that decision a was made, observation y obtained, and an a priori distribution ψ . Also, the kernel in (7.4) takes the form

$$(7.13) \quad \mathcal{K}(B \mid \psi, a) := \sum_{y \in \mathbf{Y}} V(y, \psi, a) I\{T(y, \psi, a) \in B\}, \quad B \in \mathcal{B}(\Delta),$$

while the cost \tilde{c} is computed by

$$(7.14) \quad \tilde{c}(\psi, a) := \sum_{x \in \mathbf{X}} c(x, a) \psi(x).$$

Remark 7.3. It is common to specify, instead of the kernel P , a transition kernel \bar{P} on \mathbf{X} given $\mathbf{X} \times \mathbf{A}$, and an *observation kernel* \bar{Q} on \mathbf{Y} given $\mathbf{X} \times \mathbf{A}$ [14], [63], [82], [128], [170]. Note that this is only a special case of our presentation, which happens when the kernel P admits the decomposition

$$P(x, y \mid x', a) = \bar{Q}(y \mid x, a) \bar{P}(x \mid x', a).$$

In this case, we can express (7.10)–(7.12) in a convenient vector form by viewing ψ as an element of $\mathbb{R}^{\mathbf{X}}$ and defining the transition matrix $[\bar{P}(a)]_{x, x'} := \bar{P}(x \mid x', a)$ and the observation matrix $Q_y(a) := \text{diag}\{Q(y \mid x, a) : x \in \mathbf{X}\}$. Then, with \bar{q} denoting the vector in $\mathbb{R}^{\mathbf{X}}$ defined by $\bar{q}_x(y \mid \psi, a) := q(x, y \mid \psi, a)$ and $\mathbf{1}' = (1, \dots, 1)$, we have

$$(7.10') \quad \bar{q}(y \mid \psi, a) = \bar{Q}_y(a) \bar{P}(a) \psi,$$

$$(7.11') \quad V(y, \psi, a) = \mathbf{1}' \bar{Q}_y(a) \bar{P}(a) \psi$$

(analogously for (7.12)).

Note that a *nonrandomized* separated admissible policy can be viewed as a sequence of maps $\pi_t : \Delta \rightarrow \mathbf{A}$. Then an equivalent, *completely observable*, discounted cost problem (DC') can be formulated as finding a separated admissible policy that minimizes

$$J_\beta(\psi, \pi) := E_{\psi_0}^\pi \left[\sum_{t=0}^{\infty} \beta^t \tilde{c}(\Psi_t, A_t) \right].$$

The average cost problem (AC') is analogously defined.

Note that the one-stage cost function $\tilde{c}(\psi, a)$ is linear in $\psi \in \Delta$. It is easy to show that the expectation operator corresponding to the kernel \mathcal{K} preserves concavity (convexity) [6], [50]. The following results complement those in Theorem 2.1.

THEOREM 7.1. *For a (DC') decision problem, $J_\beta^*(\cdot)$ is a concave function, for all $0 < \beta < 1$. The DCOE is given by*

$$(7.15) \quad J_\beta^*(\psi) = \min_{a \in \mathbf{A}} \left\{ \tilde{c}(\psi, a) + \beta \sum_{y \in \mathbf{Y}} V(y, \psi, a) J_\beta^*(T(y, \psi, a)) \right\},$$

and any (nonrandomized) separated stationary policy that attains the minimum above is optimal.

Remark 7.4. The optimality equation (7.15) is obtained from the general theory of CMP [15], [82]. For other results, see [5]–[7], [14], [50], [128], [161], [169], [170], [171]. Also, for a survey of relevant computational methods, see [119].

In this context, a pair (ρ, h) is said to be a solution to the ACOE if, for all $\psi \in \Delta$,

$$(7.16) \quad \rho + h(\psi) = \min_{a \in \mathbf{A}} \left\{ \tilde{c}(\psi, a) + \sum_{y \in \mathbf{Y}} V(y, \psi, a) h(T(y, \psi, a)) \right\}.$$

7.3. The vanishing discount approach. As shown in §5, for a countable state space CMP, boundedness conditions on the differential discounted value function were sufficient for solutions to the corresponding ACOE to exist. We consider here the following hypothesis.

Assumption 7.2. There exists a sequence $\beta_n \uparrow 1$, such that h_{β_n} is bounded.

Despite the fact that the model $(\Delta, \mathbf{A}, \mathcal{K}, \tilde{c})$ has a general Borel state space, it has two special features that simplify the analysis via the vanishing discount method. The first of these features is the concavity of the discounted value function, while the second is the fact that the kernel $\mathcal{K}(\cdot \mid \psi, a)$ vanishes on the complement of a countable set (for fixed ψ and a), and thus the integrals with respect to \mathcal{K} reduce to infinite sums.

For the finite state and action space case, the concavity of the discounted value function has been exploited by Platzman [136] and by Ohnishi, Mine, and Kawai [132]. These authors utilize the fact that a collection of concave functions, defined on some relatively open convex set C , which are finite and pointwise bounded, is uniformly bounded and equi-Lipschitzian relative to any closed subset of C [143, Thm. 10.6]. Thus, under Assumption 7.2, the finite dimensionality of Δ and the concavity of $h_{\beta}(\cdot)$ are used in [132], [136] to obtain a bounded solution (ρ^*, h) to the ACOE, via the vanishing discount approach. In particular, they partition Δ into its interior, its vertices, and its edges, i.e.,

$$\Delta = \bigcup_{j \in \mathcal{J}} \Delta_j.$$

Note that $|\mathcal{J}| = 2^{|\mathbf{X}|+1} - 1$ and that each set Δ_j is a relatively open convex set. Given a sequence $\beta_n \uparrow 1$, then the concavity of $h_{\beta}(\cdot)$ and Assumption 7.2 are used to obtain subsequences $\beta_n(j)$ such that $\{h_{\beta_n(j)}(\cdot)\}$ converges on Δ_j . Platzman [136] defines a metric on Δ that accomplishes this partition. Let

$$\begin{aligned} \mathcal{I}(\psi) &:= \{i \in \mathbf{X} : \psi(i) > 0\}, \quad \psi \in \Delta, \\ d(\psi_1, \psi_2) &:= 1 - \min \left\{ \frac{\psi_1(i)}{\psi_2(i)} : i \in \mathcal{I}(\psi_2) \right\}, \quad \psi_1, \psi_2 \in \Delta, \\ D(\psi_1, \psi_2) &:= \max \{d(\psi_1, \psi_2); d(\psi_2, \psi_1)\}. \end{aligned}$$

In [135, pp. 88–89], Platzman shows that $D(\cdot, \cdot)$ is a metric that leaves Δ disconnected and with components identical to the elements of the partition $\{\Delta_j\}_{j \in \mathcal{J}}$. The following is shown in [136, Lemma A.1].

LEMMA 7.1. *Let $f : \Delta \rightarrow \mathbb{R}$ be concave and bounded below; then*

$$|f(\psi_1) - f(\psi_2)| \leq \text{span}(f) D(\psi_1, \psi_2).$$

Hence, under Assumption 7.2, $\{h_{\beta}(\cdot)\}_{\beta \in (0,1)}$ is an equi-Lipschitzian family, with common Lipschitz constant given by the (smallest) uniform bound, and the Arzela–Ascoli theorem can be used as in [148] to obtain a bounded solution to the ACOE.

If the state space is infinite, the above method does not work, simply because the partition induced by the Platzman metric results in a nonseparable space. In this situation, the particular structure of the kernel has been employed in [63] to develop a theoretical framework based on the notion of *invariant* subsets (subprocesses) of a CMP, and sufficient conditions are given for the existence of solutions to the ACOE, in the case of a finite action space. The key point is to note that, if we let $B(\psi, a) := \{T(y, \psi, a) : y \in \mathbf{Y}\}$, which is a countable set since \mathbf{Y} is countable, then $\mathcal{K}(B(\psi, a) | \psi, a) = 1$. Thus, at any time $t \in \mathbb{N}_0$, the set of possible *next states* for Ψ_t is the set $\bigcup_{a \in \mathbf{A}} B(\Psi_t, a)$, which is countable, provided that \mathbf{A} is finite. This special structure has also been identified by other authors, e.g., [5, p. 187], [136, p. 369], [170, pp. 19–20].

We briefly summarize the work in [63]. The notions of *descendents*, *ancestors*, and *relatives* of a point $\psi \in \mathbf{\Delta}$ are first introduced. The descendents of ψ are defined as the smallest subset of $\mathbf{\Delta}$ containing ψ that is invariant under the action of the maps in the collection $\{T(y, \cdot, a) : y \in \mathbf{Y}, a \in \mathbf{A}\}$, while the ancestors of ψ are defined as all the points in $\mathbf{\Delta}$ that reach ψ under the application a finite sequence of these maps. Finally, the relatives of a point ψ , denoted by $\mathcal{R}_\psi^{(1)}$, is the set formed by the union of its descendents and ancestors. Note that the definition of the descendents is an extension, to the present context, of Doob's concept of *consequent sets* [45, p. 206]. Subsequently, the *genealogical tree* GT_ψ of ψ is defined by

$$GT_\psi := \bigcup_{n \in \mathbb{N}} \mathcal{R}_\psi^{(n)},$$

where the sets $\mathcal{R}_\psi^{(n)}$ are defined recursively as

$$\mathcal{R}_\psi^{(n+1)} := \bigcup_{s \in \mathcal{R}_\psi^{(n)}} \mathcal{R}_s^{(1)}, \quad n \in \mathbb{N}.$$

The descendents of a point form a countable set, but the ancestors can, in general, be uncountably many. To guarantee that the relatives and hence the genealogical tree of a point is a countable set, the following condition is introduced.

Assumption 7.3. For all $y \in \mathbf{Y}$, $a \in \mathbf{A}$, and $\psi \in \mathbf{\Delta}$, $T^{-1}(y, \psi, a)$ is a countable set.

Introduce the relation $\psi \sim \psi'$ if $GT_\psi = GT_{\psi'}$. It follows that “ \sim ” defines an *equivalence relation* on $\mathbf{\Delta}$ resulting in a partition of $\mathbf{\Delta}$ into equivalence classes that are precisely the sets GT_ψ . Under Assumptions 7.2 and 7.3, the standard diagonalization argument can be employed on each equivalence class GT_ψ to construct a pair (ρ^*, h_{GT_ψ}) that solves the ACOE on GT_ψ (the boundedness hypothesis (Assumption 7.2) can be weakened by letting the constant M depend on the equivalence class). Then, by defining $h(\psi) := h_{GT_\psi}(\psi)$ for all $\psi \in \mathbf{\Delta}$, (ρ^*, h) clearly solves the ACOE on $\mathbf{\Delta}$. One peculiarity of this approach is that the resulting function h is not guaranteed to be measurable. This is not a major problem though, since an important consequence of the particular structure (with finite action space) is that the “measurability of various objects is of no essential concern” for the equivalent problem [15, p. xi]. The approach in [63] fails when the action space \mathbf{A} is not finite.

Since the vanishing discount method relies heavily on the boundedness of the differential discounted value function, the problem of finding sufficient conditions on the cost and the kernel of the process for this to hold becomes important. Platzman [136] has given (reachability and detectability) conditions for Assumption 7.2 to hold;

however, these conditions are difficult to verify. On the other hand, many models of interest possess special properties, which allow the verification of Assumption 7.2 very easily. We examine some of these properties next.

Suppose that a partial order “ \prec_{Δ} ” has been defined on Δ and let “ $\prec_{\mathbf{A}}$ ” denote a linear order on \mathbf{A} ; we assume that \mathbf{A} is finite. We also identify \mathbf{X} with \mathbb{N}_0 and endow it with its natural ordering.

DEFINITION 7.1. Consider $((\Delta, \prec_{\Delta}), (\mathbf{A}, \prec_{\mathbf{A}}), \mathcal{K}, \tilde{c})$ and let $\psi_1, \psi_2 \in \Delta$. We state the following:

- (i) The value functions are *monotone* if

$$\psi_1 \prec_{\Delta} \psi_2 \implies J_{\beta}^*(\psi_1) \leq J_{\beta}^*(\psi_2) \quad \text{for all } 0 < \beta < 1;$$

- (ii) A (nonrandomized) stationary separated policy π is *monotone* if

$$\psi_1 \prec_{\Delta} \psi_2 \implies \pi(\psi_1) \prec_{\mathbf{A}} \pi(\psi_2).$$

Two frequently used partial orders on Δ are the *stochastic dominance* \prec_{st} and the *monotone likelihood ratio* \prec_{lr} , defined below.

DEFINITION 7.2. Let $\psi_1, \psi_2 \in \Delta$; we state the following:

- (i) $\psi_1 \prec_{st} \psi_2$ if $\sum_{i \geq q} \psi_1(i) \leq \sum_{i \geq q} \psi_2(i)$, for all $q \in \mathbf{X}$;
(ii) $\psi_1 \prec_{lr} \psi_2$ if $\psi_1(j)\psi_2(i) \leq \psi_1(i)\psi_2(j)$, for all $i, j \in \mathbf{X}$ such that $i \leq j$.

Let e^j denote the element of Δ with the j th component equal to 1, $j \in \mathbf{X}$; thus, e.g., $e^0 = (1, 0, 0, \dots)$. The following is easily shown.

LEMMA 7.2. If $\psi_1, \psi_2 \in \Delta$ and $\psi_1 \prec_{lr} \psi_2$, then $\psi_1 \prec_{st} \psi_2$. Also, for all $\psi \in \Delta$, $e^0 \prec_{lr} \psi$.

DEFINITION 7.3. An action $a_j \in \mathbf{A}$ is called a reset action if, for some $j \in \mathbf{X}$, $T(y, \psi, a_j) = e^j$, for all $y \in \mathbf{Y}$ and $\psi \in \Delta$.

A reset action a_j corresponds to the core state of the system being j , with probability one, at the next time epoch after action a_j has been taken. This type of action arises naturally in manufacturing systems subject to inspection, maintenance, and replacement. The following results derive from the work of Sondik [170]; see also [63].

LEMMA 7.3. If there exists a reset action $a_j \in \mathbf{A}$, then

$$J_{\beta}^*(\psi) - J_{\beta}^*(e^j) \leq \tilde{c}(\psi, a_j) \quad \forall \psi \in \Delta.$$

If \mathbf{X} is finite and for each $j \in \mathbf{X}$ there is a corresponding reset action, then for each $\beta \in (0, 1)$ there exists $J \in \mathbf{X}$ such that

$$0 \leq J_{\beta}^*(\psi) - J_{\beta}^*(e^J) \leq M \quad \forall \psi \in \Delta,$$

where $M := \max\{c(i, a) \mid i \in \mathbf{X}, a \in \mathbf{A}\}$.

Remark 7.5. Note that, if $J_{\beta}^*(\cdot)$ is monotone with respect to \prec_{lr} and if there is an action $a_0 \in \mathbf{A}$ that resets the state to e^0 , then $0 \leq J_{\beta}^*(\psi) - J_{\beta}^*(e^0) \leq \tilde{c}(\psi, a_0)$ uniformly in $\beta \in (0, 1)$. Furthermore, note that when \mathbf{X} is finite, a constant $M > 0$ exists such that $\tilde{c}(\psi, a_0) \leq M$, for all $\psi \in \Delta$, and thus Assumption 7.2 holds.

Models with a *replacement* action that resets the system to an “as new” state e^0 have been considered in [2], [118], [131], [132], [149], [188], [189], [191]–[195]. Related problems are those considered in [66], where a reset action to a most desirable state is available, and in [90], where (maintenance) reset actions a_j are available for all $j \neq 0$, with \mathbf{X} a finite set.

7.4. The convex analytic method. We will now briefly describe Borkar’s convex analytic approach. The action set \mathbf{A} is assumed to be any compact metric space. We also assume that c and P are continuous in a . We will consider the pathwise average cost. This cannot, in general, be written as an equivalent cost in terms of $\{\Psi_t\}$, but it is natural to propose that

$$(7.17) \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \tilde{c}(\Psi_t, A_t)$$

as a substitute. Any $\mu \in \mathcal{P}(\Delta \times \mathbf{A})$ can be decomposed as

$$(7.18) \quad \mu(d\psi, da) = \bar{\mu}(d\psi)\Phi(\psi)(da),$$

where $\bar{\mu}$ is the marginal of μ on Δ and Φ is the regular conditional law defined $\bar{\mu}$ almost surely. We always work with one arbitrary representative of this equivalence class. Define $\Gamma \subset \mathcal{P}(\Delta \times \mathbf{A})$ by

$$(7.19) \quad \begin{aligned} \Gamma &= \left\{ \mu \in \mathcal{P}(\Delta \times \mathbf{A}) \mid \begin{array}{l} \text{For } \bar{\mu}, \Phi \text{ as in (7.18), } \bar{\mu} \text{ is invariant under} \\ \text{the stationary randomized policy } \Phi \end{array} \right\} \\ &= \left\{ \mu \in \mathcal{P}(\Delta \times \mathbf{A}) \mid \int \int \int f(\psi)\mathcal{K}(d\psi \mid \psi', a)\Phi(\psi')(da)\bar{\mu}(d\psi') \right. \\ &\quad \left. = \int f d\bar{\mu} \text{ for all } f \in C_b(\Delta) \right\}. \end{aligned}$$

From (7.19) we can easily check that Γ is closed. Note that the set of invariant probability measures for the process $\{\Psi_t\}$ controlled by a stationary randomized policy Φ , when nonempty, need not be a singleton. In general, it will form a closed convex set in $\mathcal{P}(\Delta)$, the extreme points of which correspond to ergodic measures. That is, the above process with one of these extreme measures (say, μ) as the initial condition will be ergodic. Then (7.17) will almost surely equal $\int \tilde{c} d\mu$. In view of the ergodic decomposition of a stationary Markov process, this will also be the case for other invariant measures (which will be a convex combination of the ergodic ones). Define

$$\rho^* = \inf_{\mu \in \Gamma} \int \tilde{c} d\mu.$$

We assume that $\rho^* < \infty$. We consider two alternative conditions under which the above infimum will be a minimum.

Assumption 7.4 (near-monotone case). c satisfies $\lim_{i \rightarrow \infty} \inf_a c(i, a) = \infty$.

Assumption 7.5 (stable case). Assumption 5.19’ (ii) holds.

Observe that the “near-monotonicity” condition here is more restrictive than the one used in §5. We now state the following result; the proof is analogous to that of Theorem 5.10.

LEMMA 7.4. *Under either Assumption 7.4 or Assumption 7.5, the map $\mu \mapsto \int \tilde{c} d\mu$ attains its minimum on Γ .*

Define the $\mathcal{P}(\Delta \times \mathbf{A})$ -valued process $\{\eta_t\}$ by

$$\int f d\eta_t = \frac{1}{t} \sum_{m=0}^{t-1} f(\Psi_t, A_t), \quad t \geq 1, \quad f \in C_b(\Delta \times \mathbf{A}),$$

where $\{\Psi_t\}$ is governed by some policy. Again, we can prove the following analogue of Lemma 5.1.

LEMMA 7.5. *With probability 1, any limit point of $\{\eta_t\}$ in $\mathcal{P}(\Delta \times \mathbf{A})$ lies in Γ .*

Consider the near-monotone case. Suppose that, for a given sample path, a subsequence of $\{\eta_t\}$ has no limit point in $\mathcal{P}(\Delta \times \mathbf{A})$. Arguments similar to those in the proof of Theorem 5.1 can be used to show that the cost must go to $+\infty$ along this subsequence. In view of Lemma 7.5, this leads to

$$(7.20) \quad \liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \tilde{c}(\Psi_t, A_t) \geq \rho^* \quad \text{a.s.}$$

Along with Lemma 7.4, this would seem to lead to the existence of an optimal stationary randomized policy. There is, however, one catch. It is not a priori clear that any initial law for $\{\Psi_t\}$ would be in the domain of attraction of the element(s) of Γ that minimize the cost (or, for that matter, whether this domain of attraction can be reached in a finite random time from any initial law under some policy). Similar “reachability” problems surface when we try to extend the dynamic programming equations. These are circumvented under somewhat stringent conditions in [136], as we have already discussed.

Finally, we can prove the convexity of Γ . Again, it is unclear how (and whether) we can characterize the extreme points of Γ as those corresponding to stationary (nonrandomized) policies. As the ACOE is not available in this approach, the existence of an optimal stationary policy remains an open issue in general. In the stable case, it is not clear if (7.20) holds, and thus this case remains open to investigation. To sum up, the convex analytic approach to POCMP needs to be further studied.

8. Multiobjective and constrained models. An important success of the convex analytic approach discussed in §5 is in the domain of multiobjective problems, in which there is more than one cost (objective) function. We will first consider a multiobjective CMP with average cost criterion recast as a CMP with several constraints. CMP with one or multiple constraints have been studied in [1], [16], [26], [40], [41], [44], [92], [93], [97], [120], [129], [145], [146], [166]. Our presentation follows [25], [26].

We consider the case when $\mathbf{S} = \{0, 1, 2, \dots\}$; \mathbf{A} , the action space, is a prescribed compact metric space; and $P(j \mid i, a)$ is continuous in a for fixed i, j . Also, $U(i) = \mathbf{A}$ for all $i \in \mathbf{S}$. In the constrained CMP problem, we have, in addition to the cost function $c \in C_b(\mathbf{S} \times \mathbf{A})$, m additional “costs” $c_i \in C_b(\mathbf{S} \times \mathbf{A})$, $1 \leq i \leq m$ and are required to satisfy

$$(8.1) \quad a_i \leq \int c_i d\hat{\eta}(f) \leq b_i, \quad 1 \leq i \leq m$$

for prescribed numbers $b_i > a_i$, $f \in \Pi_{SD}$, and $\hat{\eta}(f) \in \mathcal{P}(\mathbf{S} \times \mathbf{A})$ is as in §5. (We are assuming all costs are bounded for simplicity. Also, we are confining our attention to Π_{SSR} ; this suffices under reasonable hypotheses, as we saw in §5.) We will assume Assumption 5.20 in §5.

Recall that $I_R = \{\hat{\eta}(f) : f \in \Pi_{SSR}\}$. Let \tilde{I}_R be the subset of I_R , where the constraints (8.1) are satisfied. Then \tilde{I}_R is closed and convex. We assume also that it is compact (this will be true under Assumption 5.20 in §5). Under this assumption, we can show, as in §5, that there exists an $f^* \in \Pi_{SR}$ that is optimal for this problem. We will now proceed to show that f^* requires randomization in at most m states.

Let $g \in C_b(\mathbf{S} \times \mathbf{A})$. For some $a \in \mathbb{R}$, let $H = I_R \cap \{\psi : \int g d\psi \leq a\}$, assumed to be nonempty. Clearly, H is closed and convex. Let $\hat{\eta}(f)$ be an extreme point of H .

Suppose that it is not an extreme point of I_R itself. Then there exist distinct measures $\hat{\eta}(f_{11}), \hat{\eta}(f_{12})$ such that at least one of them (say $\hat{\eta}(f_{11})$) is not in H , and $\hat{\eta}(f)$ is a convex combination of the two. Suppose that $\hat{\eta}(f_{21}) \in I_R \setminus H, \hat{\eta}(f_{22})$ is another such pair. Then it can be shown that $\hat{\eta}(f_{ij}), 1 \leq i, j \leq 2$ are collinear (I_R, \tilde{I}_R, H , and so on are viewed as subsets of $\mathfrak{M}(\mathcal{S} \times \mathcal{A})$, the Banach space of finite signed measures on $\mathcal{S} \times \mathcal{A}$). Therefore, all pairs of points in I_R satisfying (a) at least one of them is not in H , and (b) $\hat{\eta}(f)$ is a convex combination thereof, lie on a single straight line in $\mathfrak{M}(\mathcal{S} \times \mathcal{A})$. Let Z denote the intersection of this line with I_R . Under our hypotheses on I_R, Z is a closed finite line segment. Let $\eta(f_1), \eta(f_2)$ denote its endpoints. Then it can be shown that $\eta(f_i), i = 1, 2$ are extreme points of I_R . By Lemma 5.2, $f_i \in \Pi_{SSD}$; also, f_1 and f_2 are distinct since $\hat{\eta}(f)$ is not an extreme point of I_R . Therefore, there exists an $a' \in (0, 1)$ such that

$$\hat{\eta}(f) = a'\hat{\eta}(f_1) + (1 - a')\hat{\eta}(f_2).$$

Arguing as in the proof of Lemma 5.2, it is clear that for each $i \in \mathcal{S}$ we may take $f(i)$ to be a convex combination of $f_1(i)$ and $f_2(i)$. Let $\tilde{f} \in \Pi_{SD}$ be such that, for each $i \in \mathcal{S}, \tilde{f}(i) =$ either $f_1(i)$ or $f_2(i)$. Then, under our hypotheses (Assumption 5.20 of §5), we can show that $\hat{\eta}(\tilde{f}) \in Z$. Now consider Z as a union of two closed line segments Z_1 and Z_2, Z_1 being the line segment between $\hat{\eta}(f_1)$ and $\hat{\eta}(f)$, and Z_2 that between $\hat{\eta}(f_2)$ and $\hat{\eta}(f)$. Let $\{f'_n\}$ be a sequence in Π_{SD} , defined as follows: $f'_0 = f_1$, and

$$f'_n(i) = \begin{cases} f_2(i), & i \leq n, \\ f_1(i), & i > n. \end{cases}$$

Then, by the above considerations, $\hat{\eta}(f'_n) \in Z$. Since $f'_n \rightarrow f_2$ as $n \rightarrow \infty$, we conclude that $\hat{\eta}(f'_n) \rightarrow \hat{\eta}(f_2)$ (the map $f \mapsto \hat{\eta}(f)$ is continuous under Assumption 5.19). Thus, the sequence $\hat{\eta}(f'_n), n \geq 0$ starts in Z_1 and eventually moves into Z_2 . Let n denote the first time this happens. Then either $\hat{\eta}(f'_n) = \hat{\eta}(f)$ or $\hat{\eta}(f)$ is a convex combination of $\hat{\eta}(f'_n)$ and $\hat{\eta}(f'_{n-1})$. Since $f'_n(i) = f'_{n-1}(i)$ for $i \neq n$, the arguments employed in Lemma 5.2 show that we may take $f(i) =$ the Dirac measure at $f'_n(i)$ for $i \neq n$ and $f(n) =$ a suitable convex combination of Dirac measures at $f_1(n)$ and $f_2(n)$. We have established the following result.

THEOREM 8.1. *Each extreme point of H corresponds to an $\hat{\eta}(f)$ such that $f \in \Pi_{SR}$ satisfies the following claim: For all but at most one $i, f(i)$ is a Dirac measure at some point of \mathcal{A} . For the single remaining $i, if any, $f(i)$ is a convex combination of two such Dirac measures.$*

A variant of the above theorem leads to the following result [27].

THEOREM 8.2. *The minimum of $\nu \mapsto \int c d\nu$ on \tilde{I}_R , is attained at an $\hat{\eta}(f) \in \tilde{I}_R$, where f is either deterministic or satisfies the following claim: There are states $i_1, \dots, i_k \in \mathcal{S}$ and positive integers $n_1, \dots, n_k > 1$ such that f requires randomization among n_j values at state $i_j, 1 \leq j \leq k$; requires no randomization for the remaining states; and $\sum_{i=1}^k n_i \leq m$.*

Once this existence result is available, necessary conditions for optimality can be obtained from the standard Lagrange multiplier theory.

THEOREM 8.3. *There exist $\lambda_i, \beta_i \geq 0, 1 \leq i \leq k$ such that $\hat{\eta}(f)$, as in Theorem 8.2, minimizes*

$$\eta \mapsto F(\eta, \{\lambda_i\}, \{\beta_i\}) := \int c d\eta - \sum_{i=1}^k \lambda_i (b_i - \int c_i d\eta) - \sum_{i=1}^k \beta_i (\int c_i d\eta - a_i)$$

on I_R . Furthermore, if \tilde{I}_R has nonempty interior, the following saddle-point property holds: For all $\bar{\lambda}_i, \bar{\beta}_i \geq 0$, $1 \leq i \leq k$, $\eta \in I_R$,

$$F(\hat{\eta}(f), \{\bar{\lambda}_i\}, \{\bar{\beta}_i\}) \leq F(\hat{\eta}(f), \{\lambda_i\}, \{\beta_i\}) \leq F(\eta, \{\lambda_i\}, \{\beta_i\}).$$

Remark 8.1. The result in Theorem 8.1 cannot be improved in general. Indeed, in [26] there is a counterexample to show the nonexistence of an optimal $f \in \Pi_{SD}$ for the CMP with one constraint.

Remark 8.2. We have discussed the stable case only. Analogous results can be obtained for the near-monotone case (conditions similar to Assumption 5.18). For details, we refer to [25].

Remark 8.3. When the action set \mathbf{A} is countable, analogous results are obtained in [1].

We next consider another multiobjective CMP with AC criterion. We have m cost functions $c_i \in C_b(\mathbf{S} \times \mathbf{A})$, $1 \leq i \leq m$. All cost functions are of equal importance, and, as a result, the optimality problem cannot be recast as a constrained one. Therefore, we directly deal with the optimality problem with a vector cost criterion. This has been studied in [48], [75].

Let I_R be compact. Consider the vector cost criterion

$$\left(\int c_1 d\hat{\eta}(f), \dots, \int c_m d\hat{\eta}(f) \right), \quad \hat{\eta}(f) \in I_R.$$

In general, there need not exist an $f \in \Pi_{SSR}$ that minimizes all of $\int c_i d\hat{\eta}(f)$ over I_R . This motivates the concept of Pareto optimality. An $f \in \Pi_{SSR}$ is said to be *Pareto optimal* if there does not exist any $\bar{f} \in \Pi_{SSR}$ for which $\int c_i d\hat{\eta}(\bar{f}) \leq \int c_i d\hat{\eta}(f)$, $1 \leq i \leq m$, with inequality being strict for at least one i . Pareto optimality is clearly the minimal requirement for any reasonable notion of an optimal solution for the multiobjective problem with no priority among objectives. The Pareto optimal solutions can be characterized as follows.

THEOREM 8.4. *Any $f \in \Pi_{SSR}$ that minimizes $\sum_{i=1}^m \lambda_i \int c_i d\hat{\eta}(f)$ for some $\lambda_i > 0$, $1 \leq i \leq m$ is Pareto optimal. Conversely, any Pareto optimal $\bar{f} \in \Pi_{SSR}$ minimizes the above functional for some choice of $\lambda_i \geq 0$, $1 \leq i \leq m$.*

Remark 8.4. Note that the converse is only partial, since we have $\lambda_i \geq 0$ instead of $\lambda_i > 0$. It becomes exact if \mathbf{S} and \mathbf{A} are finite.

We often reduce a vector cost criterion as above to a scalar one by introducing a “utility function.” One such case is that of finding the “shadow minimum” for the problem of minimizing the vector cost $\nu \mapsto [\int c_1 d\nu, \dots, \int c_m d\nu] \in \mathbb{R}^m$ on I_R . Letting L denote the range of this map, L can be shown to be closed and convex. Suppose that $y_i^* = \min\{\int c_i d\nu : \nu \in I_R\}$, $1 \leq i \leq m$. Let $y^* = (y_1^*, \dots, y_m^*)$. The point y^* is called the ideal (or utopian) point. The point $x^* \in L$ that is closest to y^* is called the *shadow minimum*. This point is unique and is characterized by

$$\langle y^* - x^*, z - x^* \rangle \leq 0, \quad z \in \mathbf{S}.$$

For finite \mathbf{S} and \mathbf{A} , a combined linear-quadratic program can find x^* explicitly [75]. The point x^* is easily seen to be Pareto optimal.

9. Conclusions. We hope this paper has provided a useful presentation of the problems and techniques in average cost control of Markov processes. As is amply

clear, there is not a globally applicable approach. Instead, we expect to build a library of special tricks, a collection of simple verifiable sufficient conditions under which the problem is accessible, possibly with different techniques. Going one step further, there are the more difficult, partially observable, and multiobjective problems. Though these have seen some significant results of late, there remains much more that eludes satisfactory analysis. A similar comment applies to computational aspects and adaptive control, two topics we have not touched upon here. For computational aspects, we refer to [81], [87], [137], [180] and, for adaptive control, [26], [82], [102]. Also, we have not dealt with the vast literature on *sensitive* optimality [137], [182], nor with some other criteria, such as overtaking [111], variance sensitive [198], and weighted cost [60], [65], [99]. Finally, the discrete-time models have interesting applications to continuous-time problems, for which we refer to [14, §6.7], [109], [159], [206].

Appendix. Multifunctions and measurable selectors. Let \mathbf{V} and \mathbf{W} denote nonempty Borel spaces and let $2^{\mathbf{W}}$ denote the collection of all *nonempty* subsets of \mathbf{W} . A *multifunction* (or set-valued function) Φ from \mathbf{V} to \mathbf{W} is a map $\Phi : \mathbf{V} \rightarrow 2^{\mathbf{W}}$. The subset $\text{Dom}(\Phi) := \{v \in \mathbf{V} : \Phi(v) \neq \emptyset\}$ is called the *domain* of Φ . When $\text{Dom}(\Phi) = \mathbf{V}$, we say that the map Φ is *strict*. In what follows, we assume that Φ is a strict multifunction. If, for each $v \in \mathbf{V}$, $\Phi(v)$ is a compact (closed, measurable) subset of \mathbf{W} , then Φ is said to be *compact (closed, measurable)-valued*. A *selector* (or selection) of Φ is a function $\varphi : \mathbf{V} \rightarrow \mathbf{W}$ such that $\varphi(v) \in \Phi(v)$, for all $v \in \text{Dom}(\Phi)$. The set of (Borel) measurable selectors of Φ will be denoted by $\mathcal{S}(\Phi)$. The *graph* of Φ , denoted by $\text{Graph}(\Phi)$, is defined as

$$\text{Graph}(\Phi) := \{(v, w) : v \in \mathbf{V}, w \in \Phi(v)\}.$$

For a set $W \in 2^{\mathbf{W}}$, we define

$$\Phi^{-1}[W] := \{v \in \mathbf{V} : \Phi(v) \cap W \neq \emptyset\},$$

and we say that Φ is (Borel) *measurable* if $\Phi^{-1}[B] \in \mathcal{B}(\mathbf{V})$, for each *closed* subset B of \mathbf{W} . If Φ is *closed-valued*, then measurability of Φ implies that $\text{Graph}(\Phi) \in \mathcal{B}(\mathbf{V} \times \mathbf{W})$; furthermore, if Φ is *compact-valued*, then the converse also holds [88], [184, Thm. 4.2]. The multifunction Φ is called *upper semicontinuous* if, for every $v \in \mathbf{V}$ and every open set $G \supset \Phi(v)$, there exists a neighborhood N of v such that $\Phi(v') \subset G$, for all $v' \in N$; it is called *lower semicontinuous* if, for every $v \in \mathbf{V}$ and every open set G such that $G \cap \Phi(v) \neq \emptyset$, $\Phi^{-1}(v)$ contains an open neighborhood of v . Also, Φ is said to be *continuous* if it is both upper and lower semicontinuous.

The following result, in different variations, has been shown by several authors [15, §7.5], [47, Lemma 6, p. 38], [51, Chap. 2], [88], [154] and also summarized in [82], [184, Thm. 9.1].

THEOREM A.1. *Let Φ be a compact-valued, measurable, strict multifunction from \mathbf{V} to \mathbf{W} . Let $f : \text{Graph}(\Phi) \rightarrow \mathbb{R}$ be a measurable function, such that, for each $v \in \mathbf{V}$, $f(v, \cdot)$ is lower semicontinuous on $\Phi(v)$. Then there exists a measurable selector $\varphi^* \in \mathcal{S}(\Phi)$ such that*

$$f(v, \varphi^*(v)) = \min_{w \in \Phi(v)} \{f(v, w)\} \quad \forall v \in \mathbf{V}.$$

Let $f^ : \mathbf{V} \rightarrow \mathbb{R}$, defined by $f^*(v) := f(v, \varphi^*(v))$. If Φ is upper semicontinuous and f is bounded below, then $f^* \in \mathcal{L}(\mathbf{V})$. Also, if Φ is continuous and $f \in C_b(\mathbf{V} \times \mathbf{W})$, then $f^* \in C_b(\mathbf{V})$.*

A Tauberian theorem. The following Tauberian theorem plays a very important role in the analysis of the average cost criterion. For its proof, which is very difficult to locate in the literature in this particular format, we refer to [176].

THEOREM A.2. *Let $\{a_n\}$ be a sequence of nonnegative numbers and $\beta \in (0, 1)$. Then*

$$\begin{aligned} \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{m=0}^{N-1} a_m &\leq \liminf_{\beta \uparrow 1} (1 - \beta) \sum_{n=0}^{\infty} \beta^n a_n \\ &\leq \limsup_{\beta \uparrow 1} (1 - \beta) \sum_{n=0}^{\infty} \beta^n a_n \leq \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{m=0}^{N-1} a_m. \end{aligned}$$

Acknowledgments. The authors would like to thank the anonymous referee and Associate Editor Steven E. Shreve for their constructive criticism and comments, which helped to improve this paper. Our thanks also go to Prof. Linn I. Sennott for pointing out an error in the statement of Theorem 5.1, in an earlier draft of this manuscript. We were blessed by having an excellent typist, Joan Van Cleave, who rose above mere patience when faced with numerous revisions of our “final draft.” Finally, E. Fernández-Gaucherand wishes to thank Prof. O. Hernández-Lerma of CINVESTAV-IPN, México for useful discussions.

REFERENCES

- [1] E. ALTMAN AND A. SHWARTZ, *Markov decision problems and state-action frequencies*, SIAM J. Control Optim., 29 (1991), pp. 786–809.
- [2] V. A. ANDRIYANOV, I. A. KOGAN AND G. A. UMNNOV, *Optimal control of a partially observable discrete Markov process*, Automat. Remote Control, 4 (1980), pp. 555–561.
- [3] M. AOKI, *Optimal control of partially observable Markovian systems*, J. Franklin Inst., 280 (1965), pp. 367–386.
- [4] K. J. ARROW, T. HARRIS, AND J. MARSHAK, *Optimal inventory policy*, Econometrica, 19 (1951), pp. 250–272.
- [5] K. J. ASTRÖM, *Optimal control of Markov processes with incomplete state information*, J. Math. Anal. Appl., 10 (1965), pp. 174–205.
- [6] ———, *Optimal control of Markov processes with incomplete state information, II. The convexity of the loss function*, J. Math. Anal. Appl., 26 (1969), pp. 403–406.
- [7] ———, *Stochastic control problems*, in Mathematical Control Theory, W. A. Coppel, ed., Lecture Notes in Mathematics, Vol. 680, Springer-Verlag, Berlin, 1978, pp. 1–69.
- [8] J. A. BATHER, *Optimal decision procedures for finite Markov chains, I: Examples*, Adv. Appl. Probab., 5 (1973), pp. 328–339; II: *Communicating systems*, Adv. Appl. Probab., 5 (1973), pp. 521–540; III: *General convex systems*, Adv. Appl. Probab., 5 (1973), pp. 541–553.
- [9] R. BELLMAN, *A Markovian decision problem*, J. Math. Mech., 6 (1957), pp. 679–684.
- [10] ———, *Dynamic Programming*, Princeton University Press, Princeton, NJ, 1957.
- [11] ———, *Adaptive Control Processes: A Guided Tour*, Princeton University Press, Princeton, NJ, 1961.
- [12] R. BELLMAN AND D. BLACKWELL, *On a Particular Non-Zero Sum Game*, RM-250, RAND Corp., Santa Monica, CA, 1949.
- [13] R. BELLMAN AND J. P. LA SALLE, *On Non-Zero Sum Games and Stochastic Processes*, RM-212, RAND Corp., Santa Monica, CA, 1949.
- [14] D. P. BERTSEKAS, *Dynamic Programming: Deterministic and Stochastic Models*, Prentice-Hall, Englewood Cliffs, NJ, 1987.
- [15] D. P. BERTSEKAS AND S. E. SHREVE, *Stochastic Optimal Control: The Discrete Time Case*, Academic Press, New York, 1978.

- [16] F. J. BEUTLER AND K. W. ROSS, *Optimal policies for controlled Markov chain with a constraint*, J. Math. Anal. Appl., 112 (1985), pp. 236–256.
- [17] R. N. BHATTACHARYA AND M. MAJUMDAR, *Controlled semi-Markov model under long-run average rewards*, J. Statist. Plann. Inference, 22 (1989), pp. 223–242.
- [18] D. BLACKWELL, *Discrete dynamic programming*, Ann. Math. Statist., 33 (1962), pp. 719–726.
- [19] ———, *Discounted dynamic programming*, Ann. Math. Statist., 36 (1965), pp. 226–235.
- [20] V. S. BORKAR, *Controlled Markov chains and stochastic networks*, SIAM J. Control Optim., 21 (1983), pp. 652–666.
- [21] ———, *On minimum cost per unit time control of Markov chains*, SIAM J. Control Optim., 22 (1984), pp. 965–984.
- [22] ———, *Control of Markov chains with long-run average cost criterion*, in Stochastic Differential Systems, Stochastic Control Theory and Applications (W. Fleming and P. L. Lions, eds.), The IMA Volumes in Mathematics and Its Applications, Vol. 10, Springer-Verlag, Berlin, 1988, pp. 57–77.
- [23] ———, *A convex analytic approach to Markov decision processes*, Probab. Theory Related Fields, 78 (1988), pp. 583–602.
- [24] ———, *Control of Markov chains with long-run average cost criterion: The dynamic programming equations*, SIAM J. Control Optim., 27 (1989), pp. 642–657.
- [25] ———, *Controlled Markov chains with constraints*, Proceedings of the Workshop on Recent Advances in Modelling and Control of Stochastic Systems, Bangalore, India, January 1991, Indian Academy of Sciences, to appear.
- [26] ———, *Topics in Controlled Markov Chains*, Pitman Research Notes in Math. No. 240, Longman Scientific and Technical, Harlow, 1991.
- [27] ———, *Ergodic control of Markov chains with constraints — the general case*, preprint.
- [28] V. S. BORKAR AND M. K. GHOSH, *Ergodic and adaptive control of nearest neighbour motions*, Math. Control Signals Systems, 4 (1991), pp. 81–98.
- [29] L. D. BROWN AND R. PURVES, *Measurable selection of extrema*, Ann. Statist., 1 (1973), pp. 902–912.
- [30] R. CAVAZOS-CADENA, *Necessary and sufficient conditions for a bounded solution to the optimality equation in average reward Markov decision chains*, Systems Control Lett., 10 (1988), pp. 71–78.
- [31] ———, *Necessary conditions for the optimality equation in average reward Markov decision processes*, Appl. Math. Optim., 19 (1989), pp. 97–112.
- [32] ———, *Recent results on conditions for the existence of average optimal stationary policies*, Ann. Oper. Res., 28 (1991), pp. 3–26.
- [33] ———, *A counterexample on the optimality equation in Markov decision chains with the average cost criterion*, Systems Control Lett., 16 (1991), pp. 387–392.
- [34] R. CAVAZOS-CADENA AND L. I. SENNOTT, *Comparing recent assumptions for the existence of average optimal stationary policies*, Oper. Res. Lett., to appear.
- [35] R. YA. CHITASHVILI, *A controlled finite Markov chain with an arbitrary set of decisions*, Theory Probab. Appl., 20 (1975), pp. 839–846.
- [36] E. V. DENARDO, *A Markov decision problem*, in Mathematical Programming (T. C. Hu and S. M. Robinson, eds.), Academic Press, New York, 1973.
- [37] E. V. DENARDO AND B. L. FOX, *Multichain Markov renewal programs*, SIAM J. Appl. Math., 16 (1968), pp. 468–487.
- [38] C. DERMAN, *On sequential decisions and Markov chains*, Management Sci., 9 (1962), pp. 16–24.
- [39] ———, *Denumerable state Markov decision processes — average cost criterion*, Ann. Math. Statist., 37 (1966), pp. 1545–1553.
- [40] ———, *Finite State Markovian Decision Processes*, Academic Press, New York, 1970.
- [41] C. DERMAN AND M. KLEIN, *Some remarks on finite horizon Markovian decision models*, Oper. Res., 13 (1965), pp. 272–278.
- [42] C. DERMAN AND R. E. STRAUCH, *A note on memoryless rules for controlling sequential control processes*, Ann. Math. Statist., 37 (1966), pp. 276–278.
- [43] C. DERMAN AND A. F. VEINOTT, JR., *A solution to a countable system of equations arising in Markovian decision processes*, Ann. Math. Statist., 38 (1967), pp. 582–584.
- [44] ———, *Constrained Markov decision chains*, Management Sci., 19 (1972), pp. 389–390.
- [45] J. L. DOOB, *Stochastic Processes*, John Wiley, New York, 1953.
- [46] A. W. DRAKE, *Observation of a Markov Process Through a Noisy Channel*, Ph.D. thesis, Dept. of Electrical Engineering, MIT, Cambridge, MA, 1962.

- [47] L. DUBINS AND L. SAVAGE, *How to Gamble if You Must. Inequalities for Stochastic Processes*, McGraw-Hill, New York, 1965.
- [48] S. DURINOVIC, H. M. LEE, M. N. KATEHAKIS, AND J. A. FILAR, *Multiobjective Markov decision process with average reward criterion*, Large Scale Systems, 10 (1986), pp. 215–226.
- [49] A. DVORETZKY, J. KEIFER AND J. WOLFOWITZ, *The inventory problem*, Econometrica, 20 (1956), pp. 187–222; pp. 450–466.
- [50] E. B. DYNKIN, *Controlled random sequences*, Theory Probab. Appl., 10 (1965), pp. 1–14.
- [51] E. B. DYNKIN AND A. A. YUSHKEVICH, *Controlled Markov Processes*, Springer-Verlag, New York, 1979.
- [52] A. FEDERGRUEN, A. HORDIJK, AND H. C. TIJMS, *Recurrent conditions in denumerable state Markov decision processes*, in Dynamic Programming and Its Applications, M. L. Puterman, ed., Academic Press, New York, 1978, pp. 3–22.
- [53] ———, *Denumerable state semi-Markov decision processes with unbounded costs, average cost criterion*, Stochast. Process. Appl., 9 (1979), pp. 223–235.
- [54] A. FEDERGRUEN, P. J. SCHWEITZER, AND H. C. TIJMS, *Contraction mappings underlying undiscounted Markov decision problems*, J. Math. Anal. Appl., 65 (1978), pp. 711–730.
- [55] ———, *Denumerable undiscounted semi-Markov decision processes with unbounded rewards*, Math. Oper. Res., 8 (1983), pp. 298–313.
- [56] A. FEDERGRUEN AND H. C. TIJMS, *The optimality equation in average cost denumerable state semi-Markov decision problems, recurrence conditions and algorithms*, J. Appl. Probab., 15 (1978), pp. 356–373.
- [57] E. A. FEINBERG, *On controlled finite state Markov processes with compact control sets*, Theory Probab. Appl., 20 (1975), pp. 856–862.
- [58] ———, *The existence of a stationary ε -optimal policy for a finite Markov chain*, Theory Probab. Appl., 23 (1978), pp. 297–313.
- [59] ———, *An ε -optimal control of a finite Markov chain with an average reward criterion*, Theory Probab. Appl., 25 (1980), pp. 70–81.
- [60] ———, *Controlled Markov processes with arbitrary numerical criteria*, Theory Probab. Appl., 27 (1982), pp. 486–503.
- [61] E. FERNÁNDEZ-GAUCHERAND, *Controlled Markov Processes on the Infinite Planning Horizon: Optimal and Adaptive Control*, Ph.D. thesis, Electrical and Computer Engineering Dept., University of Texas at Austin, 1991.
- [62] E. FERNÁNDEZ-GAUCHERAND, A. ARAPOSTATHIS, AND S. I. MARCUS, *On partially observable Markov decision processes with an average cost criterion*, in Proc. 28th IEEE Conf. on Decision and Control, Tampa, FL, 1989, pp. 1267–1272.
- [63] ———, *On the average cost optimality equation and the structure of optimal policies for partially observable Markov decision processes*, Ann. Oper. Res., 29 (1991), pp. 439–470.
- [64] ———, *Remarks on the existence of solutions to the average cost optimality equation in Markov decision processes*, Systems Control Lett., 15 (1990), pp. 425–432.
- [65] E. FERNÁNDEZ-GAUCHERAND, M. K. GHOSH, AND S. I. MARCUS, *Controlled Markov processes on the infinite planning horizon with a weighted cost criterion*, Contribuciones en Probabilidad y Estadística Matemática, 3 (1992), pp. 145–162.
- [66] C. H. FINE, *A quality control model with learning effects*, Oper. Res., 36 (1988), pp. 437–444.
- [67] L. FISHER AND S. M. ROSS, *An example in denumerable decision processes*, Ann. Math. Statist., 39 (1968), pp. 674–675.
- [68] J. FLYNN, *Averaging versus discounting in dynamic programming: a counterexample*, Ann. Statist., 2 (1974), pp. 411–413.
- [69] ———, *Conditions for the equivalence of optimality criteria in dynamic programming*, Ann. Statist., 4 (1976), pp. 936–953.
- [70] ———, *On optimality criteria for dynamic programs with long finite horizons*, J. Math. Anal. Appl., 76 (1980), pp. 202–208.
- [71] N. FURUKAWA, *Markovian decision processes with compact action spaces*, Ann. Math. Statist., 43 (1972), pp. 1612–1622.
- [72] J.-P. GEORGIN, *Contrôle des chaînes de Markov sur des espaces arbitraires*, Ann. Inst. H. Poincaré Probab. Statist. Sect. B, 14 (1978), pp. 255–277.
- [73] ———, *Estimation et contrôle des chaînes de Markov sur des espaces arbitraires*, in Lecture Notes Math., 636, Springer-Verlag, Berlin, 1978, pp. 71–113.
- [74] M. K. GHOSH, *Ergodic and Adaptive Control of Markov Processes*, Ph.D. thesis, Indian Institute of Science, Bangalore, India, 1988.

- [75] M. K. GHOSH, *Markov decision processes with multiple costs*, Oper. Res. Lett., 9 (1990), pp. 257–260.
- [76] M. K. GHOSH AND S. I. MARCUS, *Ergodic control of Markov chains*, in Proc. 29th IEEE Conf. on Decision and Control, Honolulu, Hawaii, 1990, pp. 258–263.
- [77] ———, *On strong average optimality of Markov decision processes with unbounded costs*, Oper. Res. Lett., 11 (1992), pp. 99–104.
- [78] I. I. GIHMAN AND A. V. SKOROHOD, *Controlled Stochastic Processes*, Springer-Verlag, New York, 1979.
- [79] D. GILLETTE, *Stochastic games with zero stop probabilities*, Contributions to the Theory of Games, III, Annals of Math. Studies, 39, Princeton University Press, Princeton, NJ, 1957, pp. 71–187.
- [80] L. G. GUBENKO AND E. S. STATLAND, *On controlled, discrete-time Markov decision processes*, Theory Probab. Math. Statist., 7 (1975), pp. 47–61.
- [81] M. HAVIV AND M. L. PUTERMAN, *An improved algorithm for solving communicating average reward Markov decision processes*, Ann. Oper. Res., 29 (1991), pp. 229–242.
- [82] O. HERNÁNDEZ-LERMA, *Adaptive Markov Control Processes*, Springer-Verlag, New York, 1989.
- [83] ———, *Average optimality in dynamic programming on Borel spaces: Unbounded costs and controls*, preprint.
- [84] O. HERNÁNDEZ-LERMA, J. C. HENNET, AND J. B. LASSERRE, *Average cost Markov decision processes: optimality conditions*, J. Math. Anal. Appl., 158 (1991), pp. 396–406.
- [85] O. HERNÁNDEZ-LERMA AND J. B. LASSERRE, *Average cost optimal policies for Markov control processes with Borel state space and unbounded costs*, Systems Control Lett., 15 (1990), pp. 349–356.
- [86] O. HERNÁNDEZ-LERMA, R. MONTES-DE-OCA, AND R. CAVAZOS-CADENA, *Recurrence conditions for Markov decision processes with Borel state space: a survey*, Ann. Oper. Res., 29 (1991), pp. 29–46.
- [87] D. P. HEYMAN AND M. J. SOBEL, *Stochastic Models in Oper. Res., vol. II: Stochastic Optimization*, McGraw-Hill, New York, 1984.
- [88] C. J. HIMMELBERG, T. PARTHASARATHY, AND F. S. VAN VLECK, *Optimal plans for dynamic programming problems*, Math. Oper. Res., 1 (1976), pp. 390–394.
- [89] K. HINDERER, *Foundations of Non-Stationary Dynamic Programming with Discrete Time Parameters*, Lecture Notes Oper. Res. Math. Systems, Vol. 33, Springer-Verlag, Berlin, 1970.
- [90] W. J. HOPP AND S. C. WU, *Multiaction maintenance under Markovian deterioration and incomplete information*, Naval Res. Logist. Quart., 35 (1988), pp. 447–462.
- [91] A. HORDIJK, *Dynamic Programming and Markov Potential Theory*, Math. Centre Tract, No. 51, Mathematisch Centrum, Amsterdam, 1974.
- [92] A. HORDIJK AND L. C. M. KALLENBERG, *Linear programming and Markov decision chains*, Management Sci., 25 (1979), pp. 352–362.
- [93] ———, *Constrained undiscounted stochastic dynamic programming*, Math. Oper. Res., 9 (1984), pp. 276–289.
- [94] A. HORDIJK AND M. L. PUTERMAN, *On the convergence of policy iteration in undiscounted finite state Markov processes; the unichain case*, Math. Oper. Res., 12 (1987), pp. 163–176.
- [95] R. HOWARD, *Dynamic Programming and Markov Decision Processes*, MIT Press, Cambridge, MA, 1960.
- [96] G. HÜBNER, *On the fixed points of the optimal reward operator in stochastic dynamic programming with discount factor greater than one*, Zeit. Angew. Math. Mech., 57 (1977), pp. 477–480.
- [97] L. C. M. KALLENBERG, *Linear Programming and Finite Markovian Control Problems*, Math. Centre Tract, No. 148, Mathematisch Centrum, Amsterdam, 1983.
- [98] S. KARLIN, *The structure of dynamic programming models*, Naval Res. Logist. Quart., 2 (1955), pp. 285–294.
- [99] D. KRASS, J. A. FILAR, AND S. SINHA, *A weighted Markov decision process*, Oper. Res., to appear.
- [100] N. V. KRYLOV, *Construction of an optimal strategy for a finite controlled chain*, Theory Probab., 10 (1965), pp. 45–54.
- [101] P. R. KUMAR, *Simultaneous identification and adaptive control of unknown systems over finite parameter sets*, IEEE Trans. Automat. Control, AC-28 (1983), pp. 68–76.
- [102] ———, *A survey of some of results in stochastic adaptive control*, SIAM J. Control Optim., 23 (1985), pp. 329–380.

- [103] P. R. KUMAR AND P. VARAIYA, *Stochastic Systems: Estimation, Identification and Adaptive Control*, Prentice-Hall, Englewood Cliffs, NJ, 1986.
- [104] M. KURANO, *Markov decision processes with a Borel measurable cost function: the average case*, Math. Oper. Res., 11 (1986), pp. 309–320.
- [105] ———, *The existence of a minimum pair of state and policy for Markov decision processes under the hypothesis of Doeblin*, SIAM J. Control Optim., 27 (1989), pp. 296–307.
- [106] ———, *Average Cost Markov Decision Processes under the Hypothesis of Doeblin*, Report No. 9, Dept. Mathematics, Faculty of Education, Chiba, Japan, 1989.
- [107] ———, *On Optimality Inequalities in Average Cost Markov Decision Processes with Doeblin's Conditions*, Report No. 1, Dept. Mathematics, Faculty of Education, Chiba, Japan, 1990.
- [108] H. J. KUSHNER, *Stochastic Stability and Control*, Academic Press, New York, 1967.
- [109] ———, *Numerical methods for stochastic control problems in continuous time*, SIAM J. Control Optim., 28 (1990), pp. 999–1048.
- [110] B. L. LAMOND AND M. L. PUTERMAN, *Generalized inverses in discrete time Markov decision processes*, SIAM J. Math. Anal. Appl., 10 (1989), pp. 118–134.
- [111] A. LEIZAROWITZ, *Infinite horizon optimization for a finite state Markov chain*, SIAM J. Control Optim., 25 (1987), pp. 1601–1618.
- [112] G. DE LEVE, *Generalized Markov Decision Processes, Part I: Model and Method*, Math. Centre Tract, No. 3, Mathematisch Centrum, Amsterdam, 1964.
- [113] ———, *Generalized Markov Decision Processes, Part II: Probabilistic Background*, Math. Centre Tract, No. 4, Mathematisch Centrum, Amsterdam, 1964.
- [114] G. DE LEVE, A. FEDERGRUEN, AND H. C. TIJMS, *A general Markov decision method*, Adv. Appl. Probab., 9 (1977), pp. 296–335.
- [115] S. A. LIPPMAN, *Semi-Markov decision processes with unbounded rewards*, Management Sci., 19 (1973), pp. 717–731.
- [116] ———, *On dynamic programming with unbounded rewards*, Management Sci., 21 (1975), pp. 1225–1233.
- [117] M. LOËVE, *Probability Theory II*, Springer-Verlag, Berlin, 1978.
- [118] W. S. LOVEJOY, *Some monotonicity results for partially observed Markov decision processes*, Oper. Res., 35 (1987), pp. 736–743.
- [119] ———, *A survey of algorithmic methods for partially observed Markov decision processes*, Ann. Oper. Res., 28 (1991), pp. 47–66.
- [120] D.-J. MA, A. M. MAKOWSKI, AND A. SHWARTZ, *Estimation and optimal control for constrained Markov chains*, in Proc. 25th IEEE Conf. on Decision and Control, Athens, Greece, 1986, pp. 994–999.
- [121] A. MAITRA, *Dynamic Programming for Countable State Systems*, Ph.D. thesis, University of California, Berkeley, CA, 1964.
- [122] ———, *Dynamic programming for countable state systems*, Sankhyā Ser. A, 27 (1965), pp. 241–248.
- [123] ———, *Discounted dynamic programming on compact metric spaces*, Sankhyā Ser. A, 30 (1968), pp. 211–216.
- [124] P. MANDL, *Estimation and control in Markov chains*, Adv. Appl. Probab., 6 (1974), pp. 40–60.
- [125] A. MANNE, *Linear programming and sequential decisions*, Management Sci., 6 (1960), pp. 259–267.
- [126] A. MARTIN-LÖF, *Existence of a stationary control for a Markov chain maximizing the average reward*, Oper. Res., 15 (1967), pp. 866–871.
- [127] B. L. MILLER AND A. F. VEINOTT, JR., *Discrete dynamic programming with a small interest rate*, Ann. Math. Statist., 40 (1969), pp. 366–370.
- [128] G. E. MONAHAN, *A survey of partially observable Markov decision processes: theory, models, and algorithms*, Management Sci., 28 (1982), pp. 1–16.
- [129] P. NAIN AND K. W. ROSS, *Optimal priority assignment with hard constraint*, IEEE Trans. Automat. Control, AC-31 (1986), pp. 883–888.
- [130] J. NEVEU, *Mathematical Foundations of the Calculus of Probability*, Holden-Day, San Francisco, CA, 1965.
- [131] M. OHNISHI, H. KAWAI, AND H. MINE, *An optimal inspection and replacement policy under incomplete state information*, European J. Oper. Res., 27 (1986), pp. 117–128.
- [132] M. OHNISHI, H. MINE, AND H. KAWAI, *An optimal inspection and replacement policy under incomplete state information: average cost criterion*, in Stochastic Models in Reliability Theory (S. Osaki and Y. Hatoyama, eds.), Lect. Notes Econ. Math. Systems, Vol. 235, Springer-Verlag, Berlin, 1984, pp. 187–197.

- [133] S. OREY, *Limit Theorems for Markov Chain Transition Probabilities*, Van Nostrand, London, 1971.
- [134] K. R. PARTHASARATHY, *Probability Measures on Metric Spaces*, Academic Press, New York, 1967.
- [135] L. K. PLATZMAN, *Finite Memory Estimation and Control of Finite Probabilistic Systems*, Ph.D. thesis, Dept. of Electrical Engineering and Computer Science, MIT, Cambridge, MA, 1977.
- [136] ———, *Optimal infinite horizon undiscounted control of finite probabilistic systems*, SIAM J. Control Optim., 18 (1980), pp. 362–380.
- [137] M. L. PUTERMAN, *Markov decision processes*, in Handbooks in Operation Research and Management Science (D. P. Heyman and M. J. Sobel, eds.), Vol. 2, North-Holland, Amsterdam, 1990, pp. 331–434.
- [138] D. RHENIUS, *Incomplete information in Markovian decision models*, Ann. Statist., 2 (1974), pp. 1327–1334.
- [139] U. RIEDER, *Measurable selection theorems for optimization problems*, Manuscripta Math., 24 (1978), pp. 115–131.
- [140] R. K. RITT AND L. I. SENNOTT, *Optimal stationary policies in general state Markov decision chains with finite action set*, Math. Oper. Res., to appear.
- [141] D. R. ROBINSON, *Markov decision chains with unbounded costs and applications to the control of queues*, Adv. Appl. Probab., 8 (1976), pp. 159–176.
- [142] ———, *Optimality conditions for a Markov decision chain with unbounded costs*, J. Appl. Probab., 17 (1980), pp. 996–1003.
- [143] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1970.
- [144] Z. ROSBERG, P. VARAIYA, AND J. WALRAND, *Optimal control of service in tandem queues*, IEEE Trans. Automat. Control, AC-27 (1982), pp. 600–610.
- [145] K. W. ROSS, *Randomized and past dependent policies for Markov decision processes with multiple constraints*, Oper. Res., 37 (1989), pp. 474–477.
- [146] K. W. ROSS AND R. VARADARAJAN, *Markov decision processes with sample path constraints: The communicating case*, Oper. Res., 37 (1989), pp. 780–790.
- [147] S. M. ROSS, *Non-discounted denumerable Markovian decision models*, Ann. Math. Statist., 39 (1968), pp. 412–423.
- [148] ———, *Arbitrary state Markovian decision processes*, Ann. Math. Statist., 39 (1968), pp. 2118–2122.
- [149] ———, *Quality control under Markovian deterioration*, Management Sci., 17 (1971), pp. 587–596.
- [150] ———, *Introduction to Stochastic Dynamic Programming*, Academic Press, New York, 1983.
- [151] Y. SAWARAGI AND T. YOSHIKAWA, *Discrete time Markovian decision processes with incomplete state observation*, Ann. Math. Statist., 41 (1970), pp. 78–86.
- [152] M. SCHÄL, *On continuous dynamic programming with discrete time parameters*, Z. Wahrsch. Verw. Geb., 21 (1972), pp. 279–288.
- [153] ———, *On dynamic programming: compactness of the space of policies*, Stochast. Process. Appl., 3 (1975), pp. 345–364.
- [154] ———, *Conditions for optimality in dynamic programming and for the limit of n-stage optimal policies to be optimal*, Z. Wahrsch. Verw. Geb., 32 (1975), pp. 179–196.
- [155] L. I. SENNOTT, *A new condition for the existence of optimal stationary policies in average cost Markov decision processes*, Oper. Res. Lett., 5 (1986), pp. 17–23.
- [156] ———, *Average cost optimal stationary policies in infinite state Markov decision processes with unbounded costs*, Oper. Res., 37 (1989), pp. 626–633.
- [157] ———, *Average cost semi-Markov decision processes and the control of queueing systems*, Probab. Engrg. Inform. Sci., 3 (1989), pp. 247–272.
- [158] ———, *The average cost optimality equation and critical number policies*, preprint.
- [159] R. F. SERFOZO, *An equivalence between continuous and discrete time Markov decision processes*, Oper. Res., 27 (1979), pp. 616–620.
- [160] L. SHAPLEY, *Stochastic games*, Proc. Nat. Acad. Sci. U.S.A., 39 (1953), pp. 1095–1100.
- [161] A. N. SHIRYAEV, *On the theory of decision functions and control by an observation process with incomplete data*, in Selected Translations in Mathematical Statistics and Probability, Vol. 6, American Mathematical Society, Providence, RI, 1966, pp. 162–188.
- [162] ———, *Some new results in the theory of controlled random sequences*, in Selected Translations in Mathematical Statistics and Probability, Vol. 8, American Mathematical Society, Providence, RI, 1970, pp. 49–130.

- [163] A. N. SHIRYAEV, *On Markov sufficient statistics in non-additive Bayes problems of sequential analysis*, Theory Probab. Appl., 9 (1964), pp. 604–618.
- [164] S. E. SHREVE AND D. P. BERTSEKAS, *Alternative theoretical frameworks for finite horizon discrete-time stochastic optimal control*, SIAM J. Control Optim., 16 (1978), pp. 953–978.
- [165] ———, *Dynamic programming in Borel spaces*, in Dynamic Programming and Its Applications, M. L. Puterman, ed., Academic Press, New York, 1978, pp. 115–130.
- [166] A. SHWARTZ AND A. M. MAKOWSKI, *An optimal adaptive scheme for two competing queues with constraints*, in Analysis and Optimization of Systems (A. Bensoussan and J. L. Lions, eds.), Lecture Notes on Control and Information Sciences, Springer-Verlag, Berlin, 1986.
- [167] ———, *On the Poisson Equation for Markov Chains*, Report No. EE-646, Faculty of Electrical Engineering, Technion, Israel Institute of Technology, Haifa, Israel, 1987.
- [168] ———, *Comparing policies in Markov decision processes: Mandl's lemma revisited*, Math. Oper. Res., 15 (1990), pp. 155–174.
- [169] R. D. SMALLWOOD AND E. J. SONDIK, *The optimal control of partially observable Markov process over a finite horizon*, Oper. Res., 21 (1973), pp. 1071–1088.
- [170] E. J. SONDIK, *The Optimal Control of Partially Observable Markov Processes*, Ph.D. thesis, Electrical Engineering Dept., Stanford University, Stanford, CA, 1971.
- [171] ———, *The optimal control of partially observable Markov decision problems over the infinite horizon: Discounted costs*, Oper. Res., 26 (1978), pp. 282–304.
- [172] S. S. STIDHAM JR. AND R. R. WEBER, *Monotonic and insensitive optimal policies for control of queues with unbounded costs*, Oper. Res., 87 (1989), pp. 611–625.
- [173] R. E. STRAUCH, *Negative dynamic programming*, Ann. Math. Statist., 37 (1966), pp. 871–890.
- [174] C. STRIEBEL, *Sufficient statistics in the optimum control of stochastic systems*, J. Math. Anal. Appl., 12 (1965), pp. 576–593.
- [175] ———, *Optimal Control of Discrete Time Stochastic Systems*, Lecture Notes Econom. Math. Systems, Vol. 110, Springer-Verlag, Berlin, 1975.
- [176] R. SZNADJER AND J. A. FILAR, *Some comments on a theorem of Hardy and Littlewood*, J. Optim. Theory Appl., 75 (1992), to appear.
- [177] H. M. TAYLOR, *Markovian sequential replacement processes*, Ann. Math. Statist., 38 (1965), pp. 1677–1694.
- [178] L. C. THOMAS, *Connectedness conditions for denumerable state Markov decision processes*, in Recent Developments in Markov Decision Processes (R. Hartley, L. C. Thomas, and D. F. White, eds.), Academic Press, New York, 1980, pp. 181–204.
- [179] H. C. TIJMS, *On Dynamic Programming with Arbitrary State Space, Compact Action Space and the Average Reward as Criterion*, Report BW 55/75, Mathematisch Centrum, Amsterdam, 1975.
- [180] ———, *Stochastic Modelling and Analysis: A Computational Approach*, John Wiley, Chichester, UK, 1986.
- [181] J. VAN DER WAL AND J. WESSELS, *Markov decision processes*, Statist. Neerlandica, 39 (1985), pp. 219–233.
- [182] A. F. VEINOTT, *Discrete dynamic programming with sensitive discount optimality criteria*, Ann. Math. Statist., 40 (1969), pp. 1635–1660.
- [183] O. V. VISKOV AND A. N. SHIRYAEV, *On controls leading to optimal stationary models*, Trudy Mat. Inst. Steklov, 71 (1964), pp. 35–45. (In Russian.)
- [184] D. H. WAGNER, *Survey of measurable selection theorems*, SIAM J. Control Optim., 15 (1977), pp. 859–903.
- [185] H. M. WAGNER, *On the optimality of pure strategies*, Management Sci., 6 (1960), pp. 268–269.
- [186] A. WALD, *Sequential Analysis*, John Wiley, New York, 1947.
- [187] ———, *Statistical Decision Functions*, John Wiley, New York, 1950.
- [188] R. WANG, *Optimal replacement policy with unobservables states*, J. Appl. Probab., 14 (1977), pp. 340–348.
- [189] ———, *Computing optimal quality control policies — two actions*, J. Appl. Probab., 14 (1977), pp. 826–832.
- [190] R. R. WEBER AND S. S. STIDHAM JR., *Optimal control of service rates in networks of queues*, Adv. Appl. Probab., 15 (1987), pp. 202–218.
- [191] C. C. WHITE, *A Markov quality control process subject to partial observation*, Management Sci., 23 (1977), pp. 843–852.
- [192] ———, *Optimal inspection and repair of a production process subject to deterioration*, J. Oper. Res. Soc., 29 (1978), pp. 235–243.

- [193] C. C. WHITE, *Bounds on optimal cost for a replacement problem with partial observation*, Naval Res. Logist. Quart., 26 (1979), pp. 415–422.
- [194] ———, *Optimal control — limit strategies for a partially observed replacement problem*, Internat. J. Systems Sci., 10 (1979), pp. 321–331.
- [195] ———, *Monotone control laws for noisy, countable-state Markov chains*, European J. Oper. Res., 5 (1980), pp. 124–132.
- [196] C. C. WHITE AND D. J. WHITE, *Markov decision processes*, European J. Oper. Res., 39 (1989), pp. 1–16.
- [197] D. J. WHITE, *Dynamic programming of Markov chains and the method of successive approximations*, J. Math. Anal. Appl., 6 (1963), pp. 373–376.
- [198] ———, *Mean, variance, and probabilistic criteria in finite Markov decision processes: A review*, J. Optim. Theory Appl., 56 (1988), pp. 1–29.
- [199] P. WHITTLE, *Sequential decision processes with essential unobservables*, Adv. Appl. Probab., 1 (1969), pp. 271–287.
- [200] ———, *Optimization over Time: Dynamic Programming and stochastic control*, II, John Wiley, Chichester, UK, 1983.
- [201] J. WIJNGAARD, *Stationary Markovian decision problems and perturbation theory of quasicompact linear operators*, Math. Oper. Res., 2 (1977), pp. 91–102.
- [202] ———, *Existence of average optimal strategies in Markovian decision problems with strictly unbounded costs*, in Dynamic Programming and Its Applications, M. L. Puterman, ed., Academic Press, New York, 1978, pp. 369–386.
- [203] K. YAMADA, *Duality theorem in Markovian decision problems*, J. Math. Anal. Appl., 50 (1975), pp. 579–595.
- [204] A. A. YUSHKEVICH, *On a class of strategies in general Markov decision models*, Theory Probab. Appl., 18 (1973), pp. 777–779.
- [205] ———, *Reduction of a controlled Markov model with incomplete data to a problem with complete information in the case of Borel state and control spaces*, Theory Probab. Appl., 21 (1976), pp. 153–158.
- [206] ———, *On reducing a jump controllable Markov model to a model with discrete time*, Theory Probab. Appl., 25 (1980), pp. 58–59.
- [207] A. A. YUSHKEVICH AND R. YA. CHITASHVILI, *Controlled random sequences and Markov chains*, Russian Math. Surveys, 37 (1982), pp. 239–274.
- [208] H. ZIJM, *The optimality equations in multichain denumerable Markov decision processes with average cost criterion: The bounded cost case*, Statist. Decisions, 3 (1985), pp. 143–165.