

# Asymptotic Optimality and Rates of Convergence of Quantized Stationary Policies in Stochastic Control

Naci Saldi, Tamás Linder, and Serdar Yüksel

**Abstract**—We consider the discrete approximation of stationary policies for a discrete-time Markov decision process with Polish state and action spaces under total, discounted, and average cost criteria. Deterministic stationary quantizer policies are introduced and shown to be able to approximate optimal deterministic stationary policies with arbitrary precision under mild technical conditions, thus demonstrating that one can search for  $\varepsilon$ -optimal policies within the class of quantized control policies. We also derive explicit bounds on the approximation error in terms of the quantization rate.

**Index Terms**—Approximation, Markov decision processes, quantization, stationary policies, stochastic control.

## I. INTRODUCTION

In the theory of Markov decision processes (MDPs), control policies induced by measurable mappings from the state space to the action space are called stationary. For a large class of infinite horizon optimization problems, the set of stationary policies is the smallest structured set of control policies in which one can find a globally optimal policy. However, computing an optimal policy even in this class is in general computationally prohibitive for non-finite Polish (that is, complete and separable metric) state and action spaces. Furthermore, in applications to networked control, the transmission of such control actions to an actuator is not realistic when there is an information transmission constraint (imposed by the presence of a communication channel) between a plant, a controller, or an actuator.

Hence, it is of interest to study the approximation of optimal stationary policies. Several approaches have been developed in the literature to tackle this problem, most of which assume finite or countable state spaces, see [2]–[4], [17]. In this technical note, we study the following question: for infinite Borel state and action spaces, how much is lost in performance if optimal policy is represented with a finite number of bits? This formulation appears to be new in the networked control literature, where stability properties of quantized control actions have been studied extensively, but the optimization of quantized control actions has not been studied as much in the context of cost minimization.

This technical note contains two main contributions: (i) We establish conditions under which quantized control policies are asymptotically optimal; that is, as the accuracy of quantization increases, the optimal cost is achieved as the limit of the cost of quantized policies. (ii) We establish rates of convergence under further conditions; that is, we

Manuscript received April 26, 2014; revised July 14, 2014; accepted July 15, 2014. Date of publication July 28, 2014; date of current version January 21, 2015. This technical note was presented in part at the 51st Annual Allerton Conference on Communication, Control and Computing, Monticello, Illinois, October 2013. This work was supported in part by the Natural Sciences and Engineering Research Council (NSERC) of Canada. Recommended by Associate Editor P. Shi.

The authors are with the Department of Mathematics and Statistics, Queen's University, Kingston, ON, Canada (e-mail: nsaldi@mast.queensu.ca; linder@mast.queensu.ca; yuksel@mast.queensu.ca).

Digital Object Identifier 10.1109/TAC.2014.2343831

obtain bounds on the approximation loss due to quantization. These findings are somewhat analogous to results in optimal quantization theory [21].

*Organization:* In Section II we review the definition of discrete time Markov decision processes (MDP) in our setting. In Section III-A we tackle the approximation problem for the total and discounted cost cases using strategic measures. In Section III-B an analogous approximation result is obtained for the average cost case using ergodic invariant probability measures of the induced Markov chains. In Section IV we derive quantitative bounds on the approximation error in terms of the rate of the approximating quantizers for both discounted and average costs.

## II. MARKOV DECISION PROCESSES

For a metric space  $E$ , let  $\mathcal{B}(E)$  denote its Borel  $\sigma$ -algebra. Unless otherwise specified, the term “measurable” will refer to Borel measurability. We denote by  $\mathcal{P}(E)$  the set of all probability measures on  $E$ .

Consider a discrete time Markov decision process (MDP) with *state space*  $X$  and *action space*  $A$ , where  $X$  and  $A$  are complete, separable metric (Polish) spaces equipped with their Borel  $\sigma$ -algebras  $\mathcal{B}(X)$  and  $\mathcal{B}(A)$ , respectively. For all  $x \in X$ , we assume that the *set of admissible actions* is  $A$ . Let the *stochastic kernel*  $p(\cdot|x, a)$  denote the *transition probability* of the next state given that previous state-action pair is  $(x, a)$  [6]. The probability measure  $\mu$  over  $X$  denotes the initial distribution.

Define the history spaces  $H_n = (X \times A)^n \times X$ ,  $n = 0, 1, 2, \dots$  endowed with their product Borel  $\sigma$ -algebras generated by  $\mathcal{B}(X)$  and  $\mathcal{B}(A)$ . A *policy* is a sequence  $\pi = \{\pi_n\}_{n \geq 0}$  of stochastic kernels on  $A$  given  $H_n$ . A policy  $\pi$  is said to be *deterministic* if the stochastic kernels  $\pi_n$  are realized by a sequence of measurable functions  $\{f_n\}$  from  $H_n$  to  $A$ , i.e.,  $\pi_n(\cdot|h_n) = \delta_{f_n(h_n)}(\cdot)$  where  $f_n : H_n \rightarrow A$  is measurable. A policy  $\pi$  is called *stationary* if the stochastic kernels  $\pi_n$  depend only on the current state; that is,  $\pi_n = \pi_m$  ( $m, n \geq 0$ ) and  $\pi_n$  is a stochastic kernel on  $A$  given  $X$ . A policy  $\pi$  that is both deterministic and stationary is called *deterministic stationary*. Hence, deterministic stationary policies are defined by a measurable function  $f : X \rightarrow A$ . We denote by  $S$  the set of deterministic stationary policies.

According to the Ionescu Tulcea theorem [6], an initial distribution  $\mu$  on  $X$  and a policy  $\pi$  define a unique probability measure  $P_\mu^\pi$  on  $H_\infty = (X \times A)^\infty$ , which is called a *strategic measure* [5]. The expectation with respect to  $P_\mu^\pi$  is denoted by  $E_\mu^\pi$ . If  $\mu = \delta_x$  for some  $x \in X$ , we write  $P_x^\pi$  and  $E_x^\pi$  instead of  $P_{\delta_x}^\pi$  and  $E_{\delta_x}^\pi$ , respectively.

Let  $c$  and  $c_n$ ,  $n = 0, 1, 2, \dots$ , be measurable functions from  $X \times A$  to  $[0, \infty)$ . The cost functions  $w$  considered in this technical note are *expected total cost* i.e.,  $w_t(\pi, \mu) := E_\mu^\pi[\sum_{n=0}^{\infty} c_n(x_n, a_n)]$ , *expected discounted cost* i.e.,  $w_\beta(\pi, \mu) := E_\mu^\pi[\sum_{n=0}^{\infty} \beta^n c(x_n, a_n)]$  for some  $\beta \in (0, 1)$ , and *expected average cost*, i.e.,  $w_A(\pi, \mu) := \limsup_{N \rightarrow \infty} (1/N) E_\mu^\pi[\sum_{n=0}^{N-1} c(x_n, a_n)]$ . Note that the expected discounted cost is a special case of the expected total cost.

We write  $w(\pi, \mu)$  to denote the cost function (either i), ii), or iii)) of the policy  $\pi$  for the initial distribution  $\mu$ . If  $\mu = \delta_x$ , we write  $w(\pi, x)$  instead of  $w(\pi, \delta_x)$ . A policy  $\pi^*$  is called optimal if

$w(\pi^*, \mu) = \inf_{\pi} w(\pi, \mu)$  for all  $\mu \in \mathcal{P}(\mathbf{X})$ . It is well known that the set of deterministic stationary policies contains optimal policies for a large class of infinite horizon discounted cost problems (see, e.g., [6], [14]) and average cost optimal control problems (see, e.g., [1], [14]).

Throughout the technical note, the initial distribution  $\mu$  is assumed to be an arbitrary fixed distribution unless otherwise specified.

#### A. Notation and Conventions

The set of all bounded measurable real functions and bounded continuous real functions on a metric space  $\mathbf{E}$  are denoted by  $B(\mathbf{E})$  and  $C_b(\mathbf{E})$ , respectively. For any  $\nu \in \mathcal{P}(\mathbf{E})$  and measurable real function  $g$  on  $\mathbf{E}$ , define  $\nu(g) := \int g d\nu$ . Let  $\mathbf{E}_n = \prod_{i=1}^n \mathbf{E}_i$  ( $2 \leq n \leq \infty$ ) be a finite or an infinite product space. By an abuse of notation, any function  $g$  on  $\prod_{j=i_1}^{i_m} \mathbf{E}_j$ , where  $\{i_1, \dots, i_m\} \subseteq \{1, \dots, n\}$  ( $m \leq n$ ), is also treated as a function on  $\mathbf{E}_n$  by identifying it with its natural extension to  $\mathbf{E}_n$ . For any  $\pi$  and initial distribution  $\mu$ , let  $\lambda_n^{\pi, \mu}$ ,  $\lambda_{(n)}^{\pi, \mu}$ , and  $\gamma_n^{\pi, \mu}$ , respectively, denote the law of  $x_n$ ,  $(x_0, \dots, x_n)$  and  $(x_n, a_n)$  for all  $n \geq 0$ . Hence, for instance, we may write  $\lambda_{(n+1)}^{\pi, \mu}(h) = \lambda_{(n)}^{\pi, \mu}(\lambda_{(1)}^{\pi, x_n}(h))$  where  $h \in B(\mathbf{X}^{n+2})$ . Let  $\mathbb{F}$  denote the set of all measurable functions from  $\mathbf{X}$  to  $\mathbf{A}$ . For any  $g \in B(\mathbf{H}_n)$  ( $n \geq 1$ ) and  $f \in \mathbb{F}$ , define  $g_f(x_0, \dots, x_n) := g(x_0, f(x_0), \dots, f(x_{n-1}), x_n)$ . Hence, when  $c \in B(\mathbf{X} \times \mathbf{A})$ ,  $c_f(x_n) = c(x_n, f(x_n))$  since  $c \in B(\mathbf{H}_{n+1})$  by our conventions.

#### B. Problem Formulation

In this section we give a formal definition of the problems considered in this technical note. To this end, we first give the definition of a quantizer.

**Definition 2.1:** A measurable function  $q : \mathbf{X} \rightarrow \mathbf{A}$  is called a *quantizer* from  $\mathbf{X}$  to  $\mathbf{A}$  if the range of  $q$ , i.e.,  $q(\mathbf{X}) = \{q(x) \in \mathbf{A} : x \in \mathbf{X}\}$ , is finite.

The elements of  $q(\mathbf{X})$  (i.e., the possible values of  $q$ ) are called the *levels* of  $q$ . The rate  $R$  of a quantizer  $q$  is defined as the logarithm of the number of its levels:  $R = \log_2 |q(\mathbf{X})|$ . Note that  $R$  (approximately) represents the number of bits needed to losslessly encode the output levels of  $q$  using binary codewords of equal length. Let  $\mathcal{Q}$  denote the set of all quantizers from  $\mathbf{X}$  to  $\mathbf{A}$ . In this technical note we introduce a new type of policy called a *deterministic stationary quantizer policy*. Such a policy is a constant sequence  $\pi = \{\pi_n\}$  of stochastic kernels on  $\mathbf{A}$  given  $\mathbf{X}$  such that  $\pi_n(\cdot|x) = \delta_{q(x)}(\cdot)$  for all  $n$  for some  $q \in \mathcal{Q}$ . For any finite set  $\Lambda \subset \mathbf{A}$ , let  $\mathcal{Q}(\Lambda)$  denote the set of all quantizers having range  $\Lambda$  and let  $S\mathcal{Q}(\Lambda)$  denote the set of all deterministic stationary quantizer policies induced by  $\mathcal{Q}(\Lambda)$ .

The principal goal in this technical note is to determine conditions such that there exists a sequence of finite subsets  $\{\Lambda_k\}_{k \geq 1}$  of  $\mathbf{A}$  for which the following statements hold:

- (P1) For any  $\pi \in S$  there exists an approximating sequence  $\{\pi^k\}$  satisfying  $\lim_{k \rightarrow \infty} w(\pi^k, \mu) = w(\pi, \mu)$ , where  $\pi^k \in S\mathcal{Q}(\Lambda_k)$  ( $k \geq 1$ ).
- (P2) For any  $\pi \in S$  the approximating sequence  $\{\pi^k\}$  in (P1) is such that  $|w(\pi, \mu) - w(\pi^k, \mu)|$  can be explicitly upper bounded by a term depending on the cardinality of  $\Lambda_k$ .

Thus (P1) implies the existence of a sequence of stationary quantizer policies converging to an optimal stationary policy, while (P2) implies that the approximation error can be explicitly controlled.

### III. APPROXIMATION OF DETERMINISTIC STATIONARY POLICIES

A sequence  $\{\mu_n\}$  of measures on a measurable space  $(\mathbf{E}, \mathcal{E})$  is said to converge setwise [7] to a measure  $\mu$  if  $\mu_n(B) \rightarrow \mu(B)$  for all

$B \in \mathcal{E}$ , or equivalently,  $\mu_n(g) \rightarrow \mu(g)$  for all  $g \in B(\mathbf{E})$ . In this section, we will impose the following assumptions:

- (a) The stochastic kernel  $p(\cdot|x, a)$  is setwise continuous in  $a \in \mathbf{A}$ , i.e., if  $a_n \rightarrow a$ , then  $p(\cdot|x, a_n) \rightarrow p(\cdot|x, a)$  setwise for all  $x \in \mathbf{X}$ .
- (b)  $\mathbf{A}$  is compact.

**Remark 3.1:** Note that if  $\mathbf{X}$  is countable, then  $B(\mathbf{X}) = C_b(\mathbf{X})$  which implies the equivalence of setwise convergence and weak convergence. Hence, results developed in this technical note are applicable to the MDPs having weakly continuous, in the action variable, transition probabilities when the state space is countable.

**Remark 3.2:** Note that any MDP can be modeled by a discrete time dynamical system of the form  $x_{n+1} = F(x_n, a_n, v_n)$ , where the  $v_n$ 's are independent and identically distributed (i.i.d.) random variables with values in some space  $\mathbf{V}$  and common distribution  $\nu$ . In many applications, the function  $F$  has a well behaved structure and is in the form  $F(x, a, v) = H(x, a)G(v)$  or  $F(x, a, v) = H(x, a) + G(v)$ , e.g., the *fisheries management model* [6, p. 5], the *cash balance model* [15], and the *Pacific halibut fisheries management model* [18]. In these systems, assumption (a) holds for common noise processes. For instance, if  $\nu$  admits a continuous density, which is often the case in practice, then assumption (a) usually holds. We refer the reader to [15, Section 4] for a discussion on the relevance of the setwise continuity assumption on inventory control problems. In addition, the widely studied and practically important case of the additive noise system in our Example 3.1 in the next section also satisfies assumption (a).

We now define the  $ws^\infty$  topology on  $\mathcal{P}(\mathbf{H}_\infty)$  which was first introduced by Schäl in [8]. Let  $\mathcal{C}(\mathbf{H}_0) = B(\mathbf{X})$  and let  $\mathcal{C}(\mathbf{H}_n)$  ( $n \geq 1$ ) be the set of real valued functions  $g$  on  $\mathbf{H}_n$  such that  $g \in B(\mathbf{H}_n)$  and  $g(x_0, \cdot, x_1, \cdot, \dots, x_{n-1}, \cdot, x_n) \in C_b(\mathbf{A}^n)$  for all  $(x_0, \dots, x_n) \in \mathbf{X}^{n+1}$ . The  $ws^\infty$  topology on  $\mathcal{P}(\mathbf{H}_\infty)$  is defined as the smallest topology which renders all mappings  $P \mapsto P(g)$ ,  $g \in \bigcup_{n=0}^\infty \mathcal{C}(\mathbf{H}_n)$ , continuous.

Let  $d_A$  denote the metric on  $\mathbf{A}$ . Since the action space  $\mathbf{A}$  is compact and thus totally bounded, one can find a sequence of finite sets  $\{\{a_i\}_{i=1}^{m_k}\}_{k \geq 1}$  such that for all  $k$ ,  $\min_{i \in \{1, \dots, m_k\}} d_A(a, a_i) < 1/k$  for all  $a \in \mathbf{A}$ . In other words,  $\{a_i\}_{i=1}^{m_k}$  is a  $1/k$ -net in  $\mathbf{A}$ . Let  $\Lambda_k := \{a_1, \dots, a_{m_k}\}$  and for any  $f \in \mathbb{F}$  define the sequence  $\{q_k\}$  by letting

$$q_k(x) := \arg \min_{a \in \Lambda_k} d_A(f(x), a) \quad (1)$$

where ties are broken so that  $q_k$  are measurable. Note that,  $q_k \in \mathcal{Q}(\Lambda_k)$  for all  $k$  and  $q_k$  converges uniformly to  $f$  as  $k \rightarrow \infty$ . Let  $\pi \in S$  and  $\pi^k \in S\mathcal{Q}(\Lambda_k)$  be induced by  $f$  and  $q_k$ , respectively. We call each  $\pi_k$  a *quantized approximation* of  $\pi$ . In the rest of this technical note, we assume that the sequence  $\{\Lambda_k\}$ , as defined above, is fixed.

#### A. Expected Total and Discounted Costs

Recall that  $w_t$  and  $w_\beta$  denote the expected total and discounted costs, respectively. We impose the following assumptions in addition to assumptions (a) and (b):

- (c)  $c$  and  $c_n$  ( $n \geq 1$ ) are non-negative, bounded functions satisfying  $c(x, \cdot), c_n(x, \cdot) \in C_b(\mathbf{A})$  for all  $x \in \mathbf{X}$ .
- (d)  $\sup_{\tilde{\pi} \in S} \sum_{n=N+1}^\infty \gamma_n^{\tilde{\pi}, \mu}(c_n) \rightarrow 0$  as  $N \rightarrow \infty$ .

**Remark 3.3:** We note that all the results in this technical note remain valid if it is only assumed that  $c$  and  $c_n$  ( $n \geq 0$ ) are bounded and satisfies  $c(x, \cdot), c_n(x, \cdot) \in C_b(\mathbf{A})$  for all  $x \in \mathbf{X}$ .

Since the one-stage cost functions  $c_n$  are non-negative, assumption (d) is equivalent to Condition (C) in [8, pg. 349]. Clearly, the expected

discounted cost satisfies assumption (d) under assumption (c). We now state our main theorem in this subsection.

*Theorem 3.1:* Suppose assumptions (a), (b), (c) hold. Let  $\pi \in S$  and  $\{\pi^k\}$  be the quantized approximations of  $\pi$ . Then,  $w_\beta(\pi^k, \mu) \rightarrow w_\beta(\pi, \mu)$  as  $k \rightarrow \infty$ . The same statement is true for  $w_t$  if we further impose assumption (d).

The proof of Theorem 3.1 requires the following proposition which is proved in Appendix V-A.

*Proposition 3.1:* Suppose assumptions (a) and (b) hold. Then for any  $\pi \in S$ , the strategic measures  $\{P_\mu^{\pi^k}\}$  induced by the quantized approximations  $\{\pi^k\}$  of  $\pi$  converge to the strategic measure  $P_\mu^\pi$  of  $\pi$  in the  $ws^\infty$  topology. Hence,  $\gamma_n^{\pi^k, \mu}(c_n) \rightarrow \gamma_n^{\pi, \mu}(c_n)$  as  $k \rightarrow \infty$  under assumption (c).

*Proof of Theorem 3.1:* Since  $w_\beta$  is a special case of  $w_t$  and satisfies (d) under assumption (c), it is enough to prove the theorem for  $w_t$ . By Proposition 3.1,  $\gamma_n^{\pi^k, \mu}(c_n) \rightarrow \gamma_n^{\pi, \mu}(c_n)$  as  $k \rightarrow \infty$  for all  $n$ . Then, we have

$$\begin{aligned} & \limsup_{k \rightarrow \infty} |w_t(\pi^k, \mu) - w_t(\pi, \mu)| \\ & \leq \limsup_{k \rightarrow \infty} \sum_{n=0}^{\infty} |\gamma_n^{\pi^k, \mu}(c_n) - \gamma_n^{\pi, \mu}(c_n)| \\ & \leq \lim_{k \rightarrow \infty} \sum_{n=0}^N |\gamma_n^{\pi^k, \mu}(c_n) - \gamma_n^{\pi, \mu}(c_n)| + 2 \sup_{\pi \in S} \sum_{n=N+1}^{\infty} \gamma_n^{\pi, \mu}(c_n). \quad (2) \end{aligned}$$

Since the first and second terms in the last expression converge to zero as  $N \rightarrow \infty$  by Proposition 3.1 and assumption (d), respectively, the proof is complete. ■

*Remark 3.4:* Notice that this proof implicitly shows that  $w_t$  and  $w_\beta$  are sequentially continuous with respect to the strategic measures in the  $ws^\infty$  topology.

The following is a generic example frequently considered in the theory of Markov decision processes (see [12]).

*Example 3.1:* Let us consider an additive-noise system given by

$$x_{n+1} = F(x_n, a_n) + v_n, \quad n = 0, 1, 2, \dots$$

where  $X = \mathbb{R}^n$  and the  $v_n$ 's are independent and identically distributed (i.i.d.) random vectors whose common distribution has a continuous, bounded, and strictly positive probability density function. A non-degenerate Gaussian distribution satisfies this condition. We assume that the action space  $A$  is a compact subset of  $\mathbb{R}^d$  for some  $d \geq 1$ , the one stage cost functions  $c$  and  $c_n$  ( $n \geq 1$ ) satisfy assumption (c), and  $F(x, \cdot)$  is continuous for all  $x \in X$ . It is straightforward to show that assumption (a) holds under these conditions. Hence, under assumption (d) on the cost functions  $c_n$ , Theorem 3.1 holds for this system.

### B. Expected Average Cost

We are still assuming (a), (b), and (c). In contrast to the expected total and discounted cost criteria, the expected average cost is in general not sequentially continuous with respect to strategic measures for the  $ws^\infty$  topology under practical assumptions. Instead, we develop an approach based on the convergence of the sequence of invariant probability measures under quantized stationary policies.

Recall that  $w_A$  denotes the expected average cost. Observe that any deterministic stationary policy  $\pi$ , induced by  $f$ , defines a stochastic kernel on  $X$  given  $X$  via

$$Q_\pi(\cdot|x) := \lambda_1^{\pi, x}(\cdot) = p(\cdot|x, f(x)). \quad (3)$$

Let us write  $Q_\pi g(x) := \lambda_1^{\pi, x}(g)$ . If  $Q_\pi$  admits an ergodic invariant probability measure  $\nu_\pi$ , then by Theorem 2.3.4 and Proposition 2.4.2

in [7], there exists an invariant set with full  $\nu_\pi$  measure such that for all  $x$  in that set we have

$$\begin{aligned} w_A(\pi, x) &= \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} \gamma_n^{\pi, \mu}(c) \\ &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} \lambda_n^{\pi, x}(c_f) = \nu_\pi(c_f). \quad (4) \end{aligned}$$

Let  $M_\pi \in \mathcal{B}(X)$  be the set of all  $x \in X$  such that convergence in (4) holds. Hence,  $\nu_\pi(M_\pi) = 1$  if  $\nu_\pi$  exists. The following assumptions will be imposed in the main theorem of this section.

- (e) For any  $\pi \in S$ ,  $Q_\pi$  has a unique invariant probability measure  $\nu_\pi$ .
- (f1) The set  $\Gamma_S := \{\nu \in \mathcal{P}(X) : \nu Q_\pi = \nu \text{ for some } \pi \in S\}$  is relatively sequentially compact in the setwise topology.
- (f2) There exists  $x \in X$  such that for all  $B \in \mathcal{B}(X)$ ,  $\lambda_n^{\pi, x}(B) \rightarrow \nu_\pi(B)$  uniformly in  $\pi \in S$ .
- (g)  $M := \bigcap_{\pi \in S} M_\pi \neq \emptyset$ .

*Theorem 3.2:* Let the initial distribution  $\mu$  be concentrated on some  $x \in M$ . Let  $\pi \in S$  and  $\{\pi^k\}$  be the quantized approximations of  $\pi$ . Then,  $w_A(\pi^k, \mu) \rightarrow w_A(\pi, \mu)$  under the assumptions (e), (f1) or (f2), and (g).

*Proof:* See Appendix V-B. ■

In the rest of this section we will derive conditions under which assumptions (e), (f1), (f2), and (g) hold. To begin with, assumptions (e), (f2), and (g) are satisfied under any of the conditions  $Ri$ ,  $i \in \{0, 1, 1(a), 1(b), 2, \dots, 6\}$  in [10]. Moreover,  $M = X$  in (g) if at least one of the above conditions holds. The next step is to find sufficient conditions for assumptions (e), (f1), and (g) to hold.

Observe that the stochastic kernel  $p$  on  $X$  given  $X \times A$  can be written as a measurable mapping from  $X \times A$  to  $\mathcal{P}(X)$  if  $\mathcal{P}(X)$  is equipped with its Borel  $\sigma$ -algebra generated by the weak topology [7], i.e.,  $p(\cdot|x, a) : X \times A \rightarrow \mathcal{P}(X)$ . We impose the following assumption:

- (e1)  $p(\cdot|x, a) \leq \zeta(\cdot)$  for all  $x \in X$ ,  $a \in A$  for some finite measure  $\zeta$  on  $X$ .

*Proposition 3.2:* Suppose (e1) holds. Then, for any  $\pi \in S$  induced by  $f$ ,  $Q_\pi$  has an invariant probability measure  $\nu_\pi$ . Furthermore,  $\Gamma_S$  is sequentially relatively compact in the setwise topology. Hence, (e1) implies assumption (f1). In addition, if these invariant measures are unique, then assumptions (e) and (g) also hold with  $M = X$  in (g).

*Proof:* For any  $\pi \in S$ , define  $Q_{\pi, x}^{(N)}(\cdot) := (1/N) \sum_{n=0}^{N-1} \lambda_n^{\pi, x}(\cdot)$  for some  $x \in X$ . Clearly,  $Q_{\pi, x}^{(N)} \leq \zeta$  for all  $N$ . Hence, by [7, Corollary 1.4.5] there exists a subsequence  $\{Q_{\pi, x}^{(N_k)}\}$  which converges to some probability measure  $\nu_\pi$  setwise. Following the same steps in [11, Theorem 4.17] one can show that  $\nu_\pi(g) = \nu_\pi(Q_\pi g)$ , for all  $g \in \mathcal{B}(X)$ . Hence,  $\nu_\pi$  is an invariant probability measure for  $Q_\pi$ .

Furthermore, assumption (e1) implies  $\nu_\pi \leq \zeta$  for all  $\nu_\pi \in \Gamma_S$ . Thus,  $\Gamma_S$  is relatively sequentially compact in the setwise topology by again [7, Corollary 1.4.5].

Finally, for any  $\pi$ , if the invariant measure  $\nu_\pi$  is unique, then every setwise convergent subsequence of the relatively sequentially compact sequence  $\{Q_{\pi, x}^{(N)}\}$  must converge to  $\nu_\pi$ . Hence,  $Q_{\pi, x}^{(N)} \rightarrow \nu_\pi$  setwise which implies that  $w_A(\pi, x) = \limsup_{N \rightarrow \infty} Q_{\pi, x}^{(N)}(c_f) = \lim_{N \rightarrow \infty} Q_{\pi, x}^{(N)}(c_f) = \nu_\pi(c_f)$  for all  $x \in X$  since  $c_f \in \mathcal{B}(X)$ . Thus,  $M = X$  in (g). ■

*Example 3.2:* Let us consider an additive-noise system in Example 3.1 with the same assumptions. Furthermore, we assume  $F$  is bounded. Observe that for any  $\pi \in S$ , if  $Q_\pi$  has an invariant probability measure, then it has to be unique [7, Lemma 2.2.3] since there cannot exist disjoint invariant sets due to the positivity of probability density function. Since this system satisfies (e1) and  $R1(a)$  in [10]

due to the boundedness of  $F$ , assumptions (e), (f1), (f2), and (g) hold with  $\mathbf{M} = \mathbf{X}$ . This means that Theorem 3.2 holds for an additive noise system under the above conditions.

#### IV. RATES OF CONVERGENCE

Let  $\|\cdot\|_{TV}$  [7] denote the total variation distance between measures. We will impose a new set of assumptions in this section:

- (h)  $\mathbf{A}$  is an infinite compact subset of  $\mathbb{R}^d$  for some  $d \geq 1$ .
- (j)  $c$  is bounded and  $|c(x, \tilde{a}) - c(x, a)| \leq K_1 d_{\mathbf{A}}(\tilde{a}, a)$  for all  $x$ , and some  $K_1 \geq 0$ .
- (k)  $\|p(\cdot|x, \tilde{a}) - p(\cdot|x, a)\|_{TV} \leq K_2 d_{\mathbf{A}}(\tilde{a}, a)$  for all  $x$ , and some  $K_2 \geq 0$ .
- (l) There exists positive constants  $C$  and  $\kappa \in (0, 1)$  such that for all  $\pi \in S$ , there is a (necessarily unique) probability measure  $\nu_{\pi} \in \mathcal{P}(\mathbf{X})$  satisfying  $\|\lambda_n^{\pi, x} - \nu_{\pi}\|_{TV} \leq C\kappa^n$  for all  $x \in \mathbf{X}$  and  $n \geq 1$ .

Assumption (l) implies that for any policy  $\pi \in S$ , the stochastic kernel  $Q_{\pi}$ , defined in (3), has a unique invariant probability measure  $\nu_{\pi}$  and satisfies *geometric ergodicity* [7]. Note that (l) holds under any of the conditions  $Ri$ ,  $i \in \{0, 1, 1(a), 1(b), 2, \dots, 5\}$  in [10]. Moreover, one can explicitly compute the constants  $C$  and  $\kappa$  for certain systems. For instance, consider an additive-noise system in Example 3.1 with Gaussian noise. Let  $\mathbf{X} = \mathbb{R}$ . Assume  $F$  has a bounded range so that  $F(\mathbb{R}) \subset [-L, L]$  for some  $L > 0$ . Then, assumption (l) holds with  $C = 2$  and  $\kappa = 1 - \varepsilon L$ , where  $\varepsilon = (1/\sigma\sqrt{2\pi}) \exp^{-(2L)^2/2\sigma^2}$ . For further conditions that imply (l) we refer the reader to [7], [10].

The following example gives the sufficient conditions for the additive noise system under which (j), (k), and (l) hold.

*Example 4.3:* Consider the additive-noise system in Example 3.1. In addition to the assumptions there, suppose  $F(x, \cdot)$  is Lipschitz uniformly in  $x \in \mathbf{X}$  and the common density  $g$  of the  $v_n$ 's is Lipschitz on all compact subsets of  $\mathbf{X}$ . Note that a Gaussian density has these properties. Let  $c(x, a) := \|x - a\|^2$ . Under these conditions, assumptions (j) and (k) hold for the additive noise system. If we further assume that  $F$  is bounded, then assumption (l) holds as well.

The following result is a consequence of the fact that if  $\mathbf{A}$  is a compact subset of  $\mathbb{R}^d$  then there exist a constant  $\alpha > 0$  and finite subsets  $\Lambda_k \subset \mathbf{A}$  with cardinality  $|\Lambda_k| = k$  such that  $\max_{x \in \mathbf{A}} \min_{y \in \Lambda_k} d_{\mathbf{A}}(x, y) \leq \alpha(1/k)^{1/d}$  for all  $k$ , where  $d_{\mathbf{A}}$  is the Euclidean distance on  $\mathbf{A}$  inherited from  $\mathbb{R}^d$ .

*Lemma 4.1:* Let  $\mathbf{A} \subset \mathbb{R}^d$  be compact. Then for any measurable function  $f : \mathbf{X} \rightarrow \mathbf{A}$  we can construct a sequence of quantizers  $\{q_k\}$  from  $\mathbf{X}$  to  $\mathbf{A}$  which satisfy  $\sup_{x \in \mathbf{X}} d_{\mathbf{A}}(q_k(x), f(x)) \leq \alpha(1/k)^{1/d}$  for some constant  $\alpha$ .

The following proposition is the key result in this section. It is proved in Appendix V-C.

*Proposition 4.3:* Let  $\pi \in S$  and  $\{\pi^k\}$  be the quantized approximations of  $\pi$ . For any initial distribution  $\mu$  we have

$$\|\lambda_n^{\pi, \mu} - \lambda_n^{\pi^k, \mu}\|_{TV} \leq \alpha K_2 (2n - 1) \left(\frac{1}{k}\right)^{\frac{1}{d}} \quad (5)$$

for all  $n \geq 1$  under assumptions (h), (j), and (k).

##### A. Expected Discounted Cost

The proof of the following theorem essentially follows from Proposition 4.3. The proof is given in Appendix V-D.

*Theorem 4.1:* Let  $\pi \in S$  and  $\{\pi^k\}$  be the quantized approximations of  $\pi$ . For any initial distribution  $\mu$ , we have

$$|w_{\beta}(\pi, \mu) - w_{\beta}(\pi^k, \mu)| \leq K \left(\frac{1}{k}\right)^{\frac{1}{d}} \quad (6)$$

where  $K = \alpha/(1 - \beta)(K_1 - \beta K_2 M + (2\beta M K_2)/(1 - \beta))$  with  $M := \sup_{(x, a) \in \mathbf{X} \times \mathbf{A}} |c(x, a)|$  under assumptions (h), (j), and (k).

##### B. Expected Average Cost

Note that for any  $\pi \in S$ , induced by  $f$ , assumption (l) implies that  $\nu_{\pi}$  is a unique invariant probability measure for  $Q_{\pi}$  and that  $w_A(\pi, x) = \nu_{\pi}(c_f)$  for all  $x$  when  $c$  is as in the assumption (c). The following theorem basically follows from Proposition 4.3 and assumption (l). It is proved in Appendix V-E.

*Theorem 4.2:* Let  $\pi \in S$  and  $\{\pi^k\}$  be the quantized approximations of  $\pi$ . Under assumptions (h), (j), (k), and (l), for any  $x \in \mathbf{X}$  we have

$$|w_A(\pi, x) - w_A(\pi^k, x)| \leq 2MC\kappa^n + K_n \left(\frac{1}{k}\right)^{\frac{1}{d}} \quad (7)$$

for all  $n \geq 0$ , where  $K_n = ((2n - 1)K_2\alpha M + K_1\alpha)$  and  $M := \sup_{(x, a) \in \mathbf{X} \times \mathbf{A}} |c(x, a)|$ .

Observe that depending on the values of  $C$  and  $\kappa$ , we can first make the first term in the upper bound small enough by choosing sufficiently large  $n$ , and then for this  $n$  we can choose  $k$  large enough such that the second term in the upper bound is small.

*Order Optimality:* The following example demonstrates that the order of approximation errors in Theorems 4.2 and 4.1 cannot be better than  $O((1/k)^{1/d})$ . More precisely, we exhibit a simple standard example where we can lower bound the approximation errors for the optimal stationary policy by  $L(1/k)^{1/d}$ , for some positive constant  $L$ .

In what follows  $h(\cdot)$  and  $h(\cdot| \cdot)$  denote differential and conditional differential entropies, respectively [19].

*Example 4.4:* Consider the linear system

$$x_{n+1} = Ax_n + Ba_n + v_n, \quad n = 0, 1, 2, \dots$$

where  $\mathbf{X} = \mathbf{A} = \mathbb{R}^d$  and the  $v_n$ 's are i.i.d. random vectors whose common distribution has density  $g$ . For simplicity suppose that the initial distribution  $\mu$  has the same density  $g$ . It is assumed that the differential entropy  $h(g) := -\int_{\mathbf{X}} g(x) \log g(x) dx$  is finite. Let the one stage cost function be  $c(x, a) := \|x - a\|$ . Clearly, the optimal stationary policy  $\pi^*$  is induced by the identity  $f(x) = x$ , having the optimal cost  $w_i(\pi, \mu) = 0$ , where  $i \in \{\beta, A\}$ . Let  $\{\pi^k\}$  be the quantized approximations of  $\pi^*$ . Fix any  $k$  and define  $D_n := E_{\mu}^{\pi^k}[c(x_n, a_n)]$  for all  $n$ . Then, by the Shannon lower bound (SLB) [20, p. 12] we have for  $n \geq 1$

$$\begin{aligned} \log k &\geq R(D_n) \geq h(x_n) + \theta(D_n) \\ &= h(Ax_{n-1} + Ba_{n-1} + v_{n-1}) + \theta(D_n) \\ &\geq h(Ax_{n-1} + Ba_{n-1} + v_{n-1}|x_{n-1}, a_{n-1}) + \theta(D_n) \\ &= h(v_{n-1}) + \theta(D_n) \end{aligned} \quad (8)$$

where  $\theta(D_n) = -d + \log((1/(dV_d\Gamma(d)))(d/D_n)^d)$ ,  $R(D_n)$  is the rate-distortion function of  $x_n$ ,  $V_d$  is the volume of the unit sphere  $S_d = \{x : \|x\| \leq 1\}$ , and  $\Gamma$  is the gamma function. Here, (8) follows from the independence of  $v_{n-1}$  and the pair  $(x_{n-1}, a_{n-1})$ . Note that  $h(v_{n-1}) = h(g)$  for all  $n$ . Hence, we obtain  $D_n \geq L(1/k)^{1/d}$ , where  $L := (d/2)(2^{h(g)} / (dV_d\Gamma(d)))^{1/d}$ . This gives  $|w_{\beta}(\pi^*, \mu) - w_{\beta}(\pi^k, \mu)| \geq L/(1 - \beta)(1/k)^{1/d}$  and  $|w_A(\pi^*, \mu) - w_A(\pi^k, \mu)| \geq L(1/k)^{1/d}$ .

#### V. DISCUSSION

Motivated by the fact that deterministic stationary policies may not be optimal in constrained MDPs even for the discounted cost (see, e.g., [16]), these approximation results are extended to randomized policies in [22]. One direction for future work is to establish similar results for approximations where the set of admissible quantizers has a certain structure, such as the set of quantizers having convex codecells [13],

which may give rise to practical design methods. As a final remark, since setwise continuity assumption might be too restrictive in certain important cases, it is of interest to study a version of this problem where the setwise continuity assumption is replaced with the weak continuity in the state-action variables.

## APPENDIX

### A. Proof of Proposition 3.1

Suppose  $g \in \mathcal{C}(\mathsf{H}_n)$  for some  $n$ . Then we have  $P_\mu^{\pi^k}(g) = \lambda_{(n)}^{\pi^k, \mu}(g_{q_k})$  and  $P_\mu^\pi(g) = \lambda_{(n)}^{\pi, \mu}(g_f)$ . Since  $g$  is continuous in the “ $a$ ” terms by definition and  $q_k$  converges to  $f$ , we have  $g_{q_k} \rightarrow g_f$ . Hence, by [9, Theorem 2.4] it is enough to prove that  $\lambda_{(n)}^{\pi^k, \mu} \rightarrow \lambda_{(n)}^{\pi, \mu}$  setwise as  $k \rightarrow \infty$ .

We will prove this by induction. Clearly,  $\lambda_{(1)}^{\pi^k, \mu} \rightarrow \lambda_{(1)}^{\pi, \mu}$  setwise by assumption (a). Assume the claim is true for some  $n \geq 1$ . For any  $h \in B(\mathsf{X}^{n+2})$  we can write  $\lambda_{(n+1)}^{\pi^k, \mu}(h) = \lambda_{(n)}^{\pi^k, \mu}(\lambda_{(1)}^{\pi^k, x_n}(h))$  and  $\lambda_{(n+1)}^{\pi, \mu}(h) = \lambda_{(n)}^{\pi, \mu}(\lambda_{(1)}^{\pi, x_n}(h))$ . Since  $\lambda_{(1)}^{\pi^k, x_n}(h) \rightarrow \lambda_{(1)}^{\pi, x_n}(h)$  for all  $(x_0, \dots, x_n) \in \mathsf{X}^{n+1}$  by assumption (a) and  $\lambda_{(n)}^{\pi^k, \mu} \rightarrow \lambda_{(n)}^{\pi, \mu}$  setwise, we have  $\lambda_{(n+1)}^{\pi^k, \mu}(h) \rightarrow \lambda_{(n+1)}^{\pi, \mu}(h)$  by [9, Theorem 2.4] which completes the proof.

### B. Proof of Theorem 3.2

Let  $Q_\pi$  and  $Q_{\pi^k}$  be the stochastic kernels, respectively, for  $\pi$  and  $\{\pi^k\}$  defined in (3). By assumption (e),  $Q_\pi$  and  $Q_{\pi^k}$  ( $k \geq 1$ ) have unique, and so ergodic, invariant probability measures  $\nu_\pi$  and  $\nu_{\pi^k}$ , respectively. Since  $\mu$  is concentrated on some  $x \in \mathsf{M}$ , we have  $w_A(\pi^k, \mu) = \nu_{\pi^k}(c_{q_k})$  and  $w_A(\pi, \mu) = \nu_\pi(c_f)$ . Observe that  $c_{q_k}(x) \rightarrow c_f(x)$  for all  $x$  by assumption (c). Hence, if we prove  $\nu_{\pi^k} \rightarrow \nu_\pi$  setwise, then by [9, Theorem 2.4] we have  $J(\pi^k, \mu) \rightarrow J(\pi, \mu)$ . We prove this first under (f1) and then under (f2).

1) *Proof Under Assumption (f1):* We show that every setwise convergent subsequence  $\{\nu_{\pi^{k_l}}\}$  of  $\{\nu_{\pi^k}\}$  must converge to  $\nu_\pi$ . Then, since  $\Gamma_s$  is relatively sequentially compact in the setwise topology, there is at least one setwise convergent subsequence  $\{\nu_{\pi^{k_l}}\}$  of  $\{\nu_{\pi^k}\}$ , which implies the result.

Let  $\nu_{\pi^{k_l}} \rightarrow \nu$  setwise for some  $\nu \in \mathcal{P}(\mathsf{X})$ . We will show that  $\nu = \nu_\pi$  or equivalently  $\nu$  is an invariant probability measure of  $Q_\pi$ . For simplicity, we write  $\{\nu_{\pi^l}\}$  instead of  $\{\nu_{\pi^{k_l}}\}$ . Let  $g \in B(\mathsf{X})$ . Then by assumption (e) we have  $\nu_{\pi^l}(g) = \nu_{\pi^l}(Q_{\pi^l}g)$ . Since  $Q_{\pi^l}g(x) \rightarrow Q_\pi g(x)$  for all  $x$  by assumption (a) and  $\nu_{\pi^l} \rightarrow \nu$  setwise, we have  $\nu_{\pi^l}(Q_{\pi^l}g) \rightarrow \nu_\pi(Q_\pi g)$  by [9, Theorem 2.4]. On the other hand since  $\nu_{\pi^l} \rightarrow \nu$  setwise we have  $\nu_{\pi^l}(g) \rightarrow \nu(g)$ . Thus  $\nu(g) = \nu(Q_\pi g)$ . Since  $g$  is arbitrary,  $\nu$  is an invariant probability measure for  $Q_\pi$ .

2) *Proof Under Assumption (f2):* Observe that for all  $x \in \mathsf{X}$  and  $n$ ,  $\lambda_n^{\pi^k, x} \rightarrow \lambda_n^{\pi, x}$  setwise as  $k \rightarrow \infty$  since  $P_x^{\pi^k} \rightarrow P_x^\pi$  in the *ws*<sup>o</sup> topology (see Proposition 3.1). Let  $B \in \mathcal{B}(\mathsf{X})$  be given and fix some  $\varepsilon > 0$ . By assumption (f2) we can choose  $N$  large enough such that  $|\lambda_N^{\tilde{\pi}, x}(B) - \nu_{\tilde{\pi}}(B)| < \varepsilon/3$  for all  $\tilde{\pi} \in \{\pi, \pi^1, \pi^2, \dots\}$ . For this  $N$ , choose  $K$  large enough such that  $|\lambda_N^{\pi^k, x}(B) - \lambda_N^{\pi, x}(B)| < \varepsilon/3$  for all  $k \geq K$ . Thus, for all  $k \geq K$  we have  $|\nu_{\pi^k}(B) - \nu_\pi(B)| \leq |\nu_{\pi^k}(B) - \lambda_N^{\pi^k, x}(B)| + |\lambda_N^{\pi^k, x}(B) - \lambda_N^{\pi, x}(B)| + |\lambda_N^{\pi, x}(B) - \nu_\pi(B)| < \varepsilon$ . Since  $\varepsilon$  is arbitrary, we obtain  $\nu_{\pi^k}(B) \rightarrow \nu_\pi(B)$ , which completes the proof.

### C. Proof of Proposition 4.3

We will prove this result by induction. Let  $\mu$  be an arbitrary initial distribution and fix  $k$ . For  $n = 1$  the claim holds by the following

argument:

$$\begin{aligned} & \|\lambda_1^{\pi, \mu} - \lambda_1^{\pi^k, \mu}\|_{TV} \\ &= 2 \sup_{B \in \mathcal{B}(\mathsf{X})} |\mu(\lambda_1^{\pi, x}(B)) - \mu(\lambda_1^{\pi^k, x}(B))| \\ &\leq \mu(\|\lambda_1^{\pi, x} - \lambda_1^{\pi^k, x}\|_{TV}) \\ &\leq \mu(K_2 d_A(f(x), q_k(x))) \text{ (by assumption(k))} \\ &\leq \sup_{x \in \mathsf{X}} K_2 d_A(f(x), q_k(x)) \leq \left(\frac{1}{k}\right)^{\frac{1}{d}} K_2 \alpha \text{ (by Lemma 4.1).} \end{aligned}$$

Observe that the bound  $\alpha K_2 (2n - 1)(1/k)^{1/d}$  is independent of the choice of initial distribution  $\mu$  for  $n = 1$ . Assume the claim is true for  $n \geq 1$ . Then we have

$$\begin{aligned} & \|\lambda_{n+1}^{\pi, \mu} - \lambda_{n+1}^{\pi^k, \mu}\|_{TV} \\ &= 2 \sup_{B \in \mathcal{B}(\mathsf{X})} |\lambda_1^{\pi, \mu}(\lambda_n^{\pi, x_1}(B)) - \lambda_1^{\pi^k, \mu}(\lambda_n^{\pi^k, x_1}(B))| \\ &= 2 \sup_{B \in \mathcal{B}(\mathsf{X})} |\lambda_1^{\pi, \mu}(\lambda_n^{\pi, x_1}(B)) - \lambda_1^{\pi, \mu}(\lambda_n^{\pi^k, x_1}(B)) \\ &\quad + \lambda_1^{\pi, \mu}(\lambda_n^{\pi^k, x_1}(B)) - \lambda_1^{\pi^k, \mu}(\lambda_n^{\pi^k, x_1}(B))| \\ &\leq \lambda_1^{\pi, \mu}(\|\lambda_n^{\pi, x_1} - \lambda_n^{\pi^k, x_1}\|_{TV}) + 2 \|\lambda_1^{\pi, \mu} - \lambda_1^{\pi^k, \mu}\|_{TV} \quad (9) \end{aligned}$$

$$\begin{aligned} &\leq \left(\frac{1}{k}\right)^{\frac{1}{d}} (2n - 1) K_2 \alpha + 2 \left(\frac{1}{k}\right)^{\frac{1}{d}} K_2 \alpha \\ &= \alpha K_2 (2(n + 1) - 1) \left(\frac{1}{k}\right)^{\frac{1}{d}} \alpha. \quad (10) \end{aligned}$$

Here (9) follows since  $|\mu(h) - \eta(h)| \leq \|\mu - \eta\|_{TV} \sup_{x \in \mathsf{X}} |h(x)|$  and (10) follows since the bound  $\lambda K_2 (2n - 1)(1/k)^{1/d}$  is independent of the initial distribution.

### D. Proof of Theorem 4.1

For any fixed  $k$  we have

$$\begin{aligned} & |w_\beta(\pi) - w_\beta(\pi^k)| \\ &= \left| \sum_{n=0}^{\infty} \beta^n \lambda_n^{\pi, \mu}(c_f) - \sum_{n=0}^{\infty} \beta^n \lambda_n^{\pi^k, \mu}(c_{q_k}) \right| \\ &\leq \sum_{n=0}^{\infty} \beta^n \left( |\lambda_n^{\pi, \mu}(c_f) - \lambda_n^{\pi, \mu}(c_{q_k})| + |\lambda_n^{\pi, \mu}(c_{q_k}) - \lambda_n^{\pi^k, \mu}(c_{q_k})| \right) \\ &\leq \sum_{n=0}^{\infty} \beta^n \left( \sup_{x_n \in \mathsf{X}} |c_f - c_{q_k}| + \|\lambda_n^{\pi, \mu} - \lambda_n^{\pi^k, \mu}\|_{TV} M \right) \\ &\leq \sum_{n=0}^{\infty} \beta^n \left( \sup_{x_n \in \mathsf{X}} d_A(f(x_n), q_k(x_n)) K_1 + \|\lambda_n^{\pi, \mu} - \lambda_n^{\pi^k, \mu}\|_{TV} M \right) \\ &\leq \sum_{n=0}^{\infty} \beta^n \left( \left(\frac{1}{k}\right)^{\frac{1}{d}} \alpha K_1 \right) + \sum_{n=1}^{\infty} \beta^n \left( \left(\frac{1}{k}\right)^{\frac{1}{d}} (2n - 1) K_2 \alpha M \right) \\ &= \left(\frac{1}{k}\right)^{\frac{1}{d}} \frac{\alpha}{1 - \beta} \left( K_1 - \beta K_2 M + \frac{2\beta M K_2}{1 - \beta} \right). \quad (11) \end{aligned}$$

Here (11) follows from Assumption (j), Proposition 4.3, and Lemma 4.1, completing the proof.

### E. Proof of Theorem 4.2

For any  $k$  and  $x \in \mathbb{X}$ , we have

$$\begin{aligned}
& |w_A(\pi, x) - w_A(\pi^k, x)| = |\nu_\pi(c_f) - \nu_{\pi^k}(c_{q_k})| \\
& \leq |\nu_\pi(c_f) - \nu_\pi(c_{q_k})| + |\nu_\pi(c_{q_k}) - \nu_{\pi^k}(c_{q_k})| \\
& \leq \sup_{x \in \mathbb{X}} K_1 d_A(f(x), q_k(x)) + \|\nu_\pi - \nu_{\pi^k}\|_{TV} M \quad (\text{by (j)}) \\
& \leq \left(\frac{1}{k}\right)^{\frac{1}{d}} K_1 \alpha + \left( \|\nu_\pi - \lambda_n^{\pi, x}\|_{TV} + \|\lambda_n^{\pi, x} - \lambda_n^{\pi^k, x}\|_{TV} \right. \\
& \quad \left. + \|\lambda_n^{\pi^k, x} - \nu_{\pi^k}\|_{TV} \right) M \\
& \leq \left(\frac{1}{k}\right)^{\frac{1}{d}} K_1 \alpha + \left( 2C\kappa^n + \left(\frac{1}{k}\right)^{\frac{1}{d}} (2n-1)K_2 \alpha \right) M \\
& = 2MC\kappa^n + ((2n-1)K_2 \alpha M + K_1 \alpha) \left(\frac{1}{k}\right)^{\frac{1}{d}} \tag{12}
\end{aligned}$$

where (12) follows from assumption (I) and Proposition 4.3.

### REFERENCES

- [1] V. Borkar, "Convex analytic methods in Markov decision processes," in *Handbook of Markov Decision Processes*, E. Feinberg and A. Shwartz, Eds. Norwell, MA: Kluwer Academic Publisher, 2002.
- [2] D. Bertsekas and J. Tsitsiklis, *Neuro-Dynamic Programming*. Boston, MA: Athena Scientific, 1996.
- [3] R. Ortner, "Pseudometrics for state aggregation in average reward Markov decision processes," in *Algorithmic Learning Theory*. Berlin, Germany: Springer-Verlag, 2007.
- [4] R. Cavazos-Cadena, "Finite-state approximations for denumerable state discounted Markov decision processes," *Appl. Math. Optim.*, vol. 14, pp. 1–26, 1986.
- [5] E. Feinberg, "On measurability and representation of strategic measures in Markov decision processes," *Stat., Prob. Game Theory*, vol. 30, pp. 29–43, pp. 29–43, 1996.
- [6] O. Hernández-Lerma and J. Lasserre, *Discrete-Time Markov Control Processes: Basic Optimality Criteria*. New York: Springer, 1996.
- [7] O. Hernández-Lerma and J. Lasserre, *Markov Chains and Invariant Probabilities*. Boston, MA: Birkhäuser, 2003.
- [8] M. Schäl, "On dynamic programming: Compactness of the space of policies," *Stochastic Process. Appl.*, vol. 3, no. 4, pp. 345–364, 1975.
- [9] R. Serfozo, "Convergence of Lebesgue integrals with varying measures," *Sankhya Ser. A*, pp. 380–402, 1982.
- [10] O. Hernández-Lerma, R. Montes-De-Oca, and R. Cavazos-Cadena, "Recurrence conditions for Markov decision processes with Borel state space: A survey," *Annu. Oper. Res.*, vol. 28, no. 1, pp. 29–46, 1991.
- [11] M. Hairer, "Ergodic properties of Markov processes," *Lecture Notes*, 2006.
- [12] O. Hernández-Lerma and R. Romera, "Limiting discounted-cost control of partially observable stochastic systems," *SIAM J. Control Optim.*, vol. 40, no. 2, pp. 348–369, 2001.
- [13] A. György and T. Linder, "Codecell convexity in optimal entropy-constrained vector quantization," *IEEE Trans. Inform. Theory*, vol. 49, no. 7, pp. 1821–1828, Jul. 2003.
- [14] E. A. Feinberg, P. O. Kasyanov, and N. V. Zadoianchuk, "Average cost Markov decision processes with weakly continuous transition probabilities," *Math. Oper. Res.*, vol. 37, no. 4, pp. 591–607, Nov. 2012.
- [15] E. A. Feinberg and M. E. Lewis, "Optimality inequalities for average cost Markov decision processes and the stochastic cash balance problem," *Math. Oper. Res.*, vol. 32, no. 4, pp. 769–783, Nov. 2007.
- [16] A. B. Piunovskiy, *Optimal Control of Random Sequences in Problems With Constraints*. New York: Kluwer, 1997.
- [17] F. Dufour and T. Prieto-Rumeau, "Finite linear programming approximations of constrained discounted Markov decision processes," *SIAM J. Control Optim.*, vol. 51, no. 2, pp. 1298–1324, 2013.
- [18] F. Dufour and T. Prieto-Rumeau, "Stochastic approximations of constrained discounted Markov decision processes," *J. Math. Anal. Appl.*, vol. 413, pp. 856–879, 2014.
- [19] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 2006.
- [20] Y. Yamada, S. Tazaki, and R. M. Gray, "Asymptotic performance of block quantizers with difference distortion measures," *IEEE Trans. Inform. Theory*, vol. 26, pp. 6–14, Jan. 1980.
- [21] R. M. Gray and D. L. Neuhoff, "Quantization," *IEEE Trans. Inform. Theory*, vol. 44, no. 6, pp. 2325–2383, Oct. 1998.
- [22] N. Saldi, T. Linder, and S. Yüksel, Quantized Stationary Control Policies in Markov Decision Processes, 2014, arXiv:1310.5770.