# The Projected Bellman Equation in Reinforcement Learning

Sean Meyn

Abstract-Q-learning has become an important part of the reinforcement learning toolkit since its introduction in the dissertation of Chris Watkins in the 1980s. In the original tabular formulation, the goal is to compute exactly a solution to the discounted-cost optimality equation, and thereby obtain the optimal policy for a Markov Decision Process. The goal today is more modest: obtain an approximate solution within a prescribed function class.

The standard algorithms are based on the same architecture as formulated in the 1980s, with the goal of finding a value function approximation that solves the so-called projected Bellman equation. While reinforcement learning has been an active research area for over four decades, there is little theory providing conditions for convergence of these Q-learning algorithms, or even existence of a solution to this equation.

The purpose of this paper is to show that a solution to the projected Bellman equation does exist, provided the function class is linear and the input used for training is a form of  $\varepsilon$ -greedy policy with sufficiently small  $\varepsilon$ . Moreover, under these conditions it is shown that the Q-learning algorithm is stable, in terms of bounded parameter estimates. Convergence remains one of many open topics for research.

#### I. Introduction

Much of reinforcement learning concerns optimal control of state space models, typically cast in a Markov Decision Process (MDP) setting. Following standard notation from the control systems literature, the state process is denoted X = $\{X_k: k \geq 0\}$ , the input process  $U = \{U_k: k \geq 0\}$ , and  $c(X_k, U_k)$  denotes the one-stage cost at time k.

This paper concerns Q-learning algorithms, motivated by the same objective as in the first formulation of Watkins [54], [53]: the infinite-horizon optimal control problem, with stateaction value function

$$Q^{\star}(x,u) = \min \sum_{k=0}^{\infty} \gamma^{k} \mathsf{E}[c(X_{k}, U_{k}) \mid X_{0} = x \,, \ U_{0} = u] \tag{1}$$

where  $\gamma \in (0,1)$  is the discount factor. The minimum in (1) is over all history dependent input sequences. This is the Ofunction of Q-learning.

Under standard assumptions an optimal input is obtained by state feedback,  $U_k^* = \phi^*(X_k^*)$  for each k, where an optimal policy  $\phi^* : X \to U$  is obtained via  $\phi^*(x) \in \arg\min_u Q^*(x, u)$ for each x [8]. To avoid technicalities (in particular to avoid discussion of measurability), it is assumed in this paper that

Financial support from ARO award W911NF2010055 and NSF award CCF 2306023 is gratefully acknowledged. Many thanks to Caio Lauand at UF for comments on a draft manuscript.

S. Meyn is with the Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL 32611, USA (email: meyn@ece.ufl.edu).

the state process evolves on a finite state space X, and the input takes values a finite set U.

The Q-function solves the Bellman equation  $Q^* = \mathcal{T}Q^*$ , in which the *Bellman operator*  $\mathcal{T}$  acts on functions  $H: X \times U \rightarrow$ 

$$TH(x, u) = c(x, u) + \mathsf{E}[\underline{H}(X_{n+1}) \mid X_n = x, \ U_n = u]$$

where throughout the paper  $H(x) := \min_{u} H(x, u), x \in X$ . It is helpful to express the Bellman equation in sample path form: for any adapted input, and k > 0,

$$Q^{\star}(X_k, U_k) = c(X_k, U_k) + \gamma \mathsf{E}[Q^{\star}(X_{k+1}) \mid \mathcal{F}_k] \tag{2}$$

in which  $\mathcal{F}_k = \sigma\{X_i, U_i : i \leq k\}$  is the history up to time k. The objective of Q-learning is to obtain an approximate solution among a parameterized class  $\{Q^{\theta}: \theta \in \mathbb{R}^d\}$ . Given

an approximation we obtain a policy defined in analogy with the optimal policy,

$$\Phi^{\theta}(x) \in \arg\min_{u} Q^{\theta}(x, u), \quad x \in \mathsf{X}, \tag{3}$$

with some fixed rule in place in case of ties.

The most common criterion for success is the solution to the projected Bellman equation: find  $\theta^* \in \mathbb{R}^d$  such that

$$0 = \mathsf{E} \big[ \{ c(X_n, U_n) + \gamma \underline{Q}^{\theta^*}(X_{n+1}) - Q^{\theta^*}(X_n, U_n) \} \zeta_n^{\theta^*} \big] \ \, (4)$$

in which the expectation is in steady-state, and  $\zeta_n^{\theta} =$  $\nabla_{\theta}Q^{\theta}(X_n,U_n)$  for each n and  $\theta \in \mathbb{R}^d$ . Alternatives are discussed in Section IV-D.

The theoretical results in this paper are obtained in the special case of linear function approximation:

$$Q^{\theta} = \theta^{\mathsf{T}} \psi \qquad \text{giving} \quad \zeta_n^{\theta} = \psi(X_n, U_n) \,,$$
 (5)

with  $\psi$  a vector of basis functions. In this case the projected Bellman equation may be expressed in the Hilbert space notation of [50],

$$Q^{\theta^*} = \Pi \mathcal{T} Q^{\theta^*}$$

in which  $\Pi$  denotes the projection onto the d-dimensional subspace  $L_2(\pi_{\theta^*})$ ; the definition of the probability mass function (pmf)  $\pi_{\theta^*}$  may be found below (39b).

Much of the present article focuses on a generalization of the original algorithm of Watkins: For initialization  $\theta_0 \in \mathbb{R}^d$ , define the sequence of estimates recursively:

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \mathcal{D}_{n+1} \zeta_n \tag{6a}$$

$$\mathcal{D}_{n+1} = c(X_n, U_n) + \gamma Q^{\theta_n}(X_{n+1}) - Q^{\theta_n}(X_n, U_n), \quad (6b)$$

in which  $\{\alpha_n\}$  is a non-negative step-size sequence,  $\{\zeta_n :=$  $\zeta_n^{\theta_n}$  are known as the *eligibility vectors* (entirely analogous to the eligibility vectors used in the TD(0) algorithm [46], [50]), and  $\{\mathcal{D}_{n+1}\}$  is known as the temporal difference sequence.

The recursion (6a) reduces to Watkins' algorithm when using a tabular basis [54], [53] (see Section III-A for definitions). See [45], [47], [37] for a range of interpretations of the algorithm.

Soon after Q-learning was introduced, it was recognized that the algorithm can be cast within the framework of stochastic approximation (SA) [49], [24]. To explain the contributions and approach to analysis in this paper it is necessary to first explain why (6a) can also be cast as an SA recursion, subject to mild assumptions on the input used for training.

**Some history** The central open issue motivating the research surveyed in this paper is this: *it is not known if the projected Bellman equation* (4) *has a solution outside of very special cases.* 

Success stories surveyed in [47] include the special case of binning [24], which is a generalization of the tabular setting, and the criterion in [35] and its improvement in [28], for which the assumptions are not easily verified in practice. The progress report in [47, Section 3.3.2] states that the *only known convergence result is due to Melo et al.* [35]. See [45, Section 11.2] for further discussion, and [22] for recent insight.

This open problem was a topic of discussion throughout the Simons program on reinforcement learning held in 2020, especially during the bootcamp lectures [48].

Thms. IV.1 and IV.5 resolve this open problem for Q-learning with optimistic training. Following several preliminaries, the proof of Thm. IV.1 is similar to the proof of convergence of  $TD(\lambda)$  learning from the dissertation of Van Roy [50], [51], and the assumptions are related to the assumptions in this prior work, even though the setting is very different.

An approximate projected Bellman equation is considered in [17], in which the minimum defining  $\underline{Q}^{\theta_n}$  is replaced with a soft-min. Under assumptions similar to those imposed here they establish the existence of a solution [17, Theorem 5.1]. This result is similar to Prop. IV.2 (ii) of the present paper.

The recent paper [29] considers Q-learning with linear function approximation and oblivious training (meaning that the input used for training does not depend directly on parameter estimates). With sufficiently large regularization they obtain a unique equilibrium for the algorithm that approximates the solution to the projected Bellman equation.

Also recent is the work of [15], which is cast in a similar setting: Q-learning with linear function approximation and oblivious training. It is argued that the use of a target network combined with a carefully constructed projection of parameters improves performance, and their error bounds are consistent with their claims. While the paper is a significant step forward, they leave open the question of existence of a solution to the projected Bellman equation. With vanishing step-size, if convergence is established with or without a target network, the limit must be a solution to the projected Bellman equation (see [37, Proposition 5.10] for proof in the case of deterministic optimal control—the arguments in the stochastic setting are identical).

The lack of theory motivated Baird's gradient descent approach [4] as well as GQ learning [31], in which the root

finding problem is replaced with the minimization of a loss function. See [3] for recent theory and Section IV-D for further discussion.

Zap stochastic approximation was introduced to ensure convergence, and also provide acceleration [20]. While originally proposed for Q-learning with linear function approximation, it was later shown to be convergent even with nonlinear function approximation [14], and the general technique applies to any application in which stochastic approximation is used. The Zap-Zero algorithm introduced in [37] and improved recently in [38] is designed to avoid matrix inversion.

Much recent research has focused on *linear* MDPs, notably [55], [25], in which the system dynamics are partially known: for a known "feature map"  $\phi \colon \mathsf{X} \times \mathsf{U} \to \mathbb{R}^d$  and an unknown sequence of probability measures  $\{\mu_i : 1 \le i \le d\}$  on  $\mathsf{X}$ , a linear MDP is assumed to have a controlled transition matrix of the form  $P_u(x,x') = \sum \phi_i(x,u)\mu_i(x')$ . There is now a relatively complete theory for this special case, in which the algorithm is designed based on knowledge of the feature map.

The reader is encouraged to see [5], [30], [33], [32] for new approaches to Q-learning based on convex programming approaches to MDPs. It is hoped that the analytical techniques presented in this paper may be adapted to these new algorithms.

**Overview** Section II surveys relevant recent results from stochastic approximation theory, and Section III provides a review of Q-learning, for which the vast majority of theory is restricted to linear function approximation and oblivious training.

Consideration of optimistic policies is postponed to Section IV, which contains the main contributions of the paper: if a smooth approximation of the  $\varepsilon$ -greedy policy is used for training, then under mild conditions the parameter estimates are bounded, and there exists a solution to the projected Bellman equation (see Thm. IV.1).

#### II. BACKGROUND AND ASSUMPTIONS

This section is devoted to three topics: assumptions surrounding the MDP model, a brief summary of results from the theory of stochastic approximation, followed by assumptions surrounding the Q-learning algorithms to be considered.

**Notation:**  $Q^*$ : Discounted-cost value function, (1).

- $f_{n+1}$ ,  $c_n$  and  $\psi_{(n)}$ , (38).
- $\phi^{\theta}$ :  $Q^{\theta}$ -greedy policy, (3).  $\widetilde{\phi}^{\theta}$  randomized policy, (31).
- $\mathcal{C}^{\Theta}$ : region of policy continuity, (29).
- $H(x) := \min_{u} H(x, u)$ , appearing in DCOE (2).
- Q-learning notation from (6b): step-size  $\alpha_n$ , eligibility vector  $\zeta_n$ , temporal difference  $\mathcal{D}_{n+1}$ .
- $\theta^* \in \mathbb{R}^d$ : solves projected Bellman equation (4).
- $\theta_n$  parameter estimate,  $\theta_n^{PR}$  PR-average, (12).
- Errors:  $\tilde{\theta}_n = \theta_n \theta^*$ ,  $\tilde{\theta}_n^{PR} = \theta_n^{PR} \theta^*$ .
- $P_u$ : controlled transition matrix, (8).
- $\bar{f}: \mathbb{R}^d \to \mathbb{R}^d$ , vector field for mean-flow, (10b).
- $\bar{f}_{\infty}$ , vector field for ODE@ $\infty$ , (22).

- $\vartheta_t$ : solution to mean-flow, (11).
- $A = \partial_{\theta} \overline{f}$  and  $A^* := A(\theta^*)$ , (19).
- $\Sigma_{\Theta}$  asymptotic covariance, (24a).
- $\Sigma_{\Theta}^{PR} = G \Sigma_{\Delta}^* G^{\mathsf{T}}$  with  $G = -(A^*)^{-1}$  and  $\Sigma_{\Delta}^*$ , (24b).
- $U_k = (1 B_k)\mathcal{U}_k + B_k\mathcal{W}_k$  training policy, (30).

## A. Markov Decision Process

While the search for an optimal policy may be restricted to static state feedback under the assumptions imposed below, in reinforcement learning it is standard practice to introduce randomization in policies as a way of introducing exploration during training. We restrict to randomized policies of the form,

$$U_k = \phi(X_k, \theta_k, I_k), \qquad k \ge 0, \tag{7}$$

in which  $I = \{I_1, I_2, ...\}$  is an i.i.d. sequence. Under the assumption that X and U are finite, we can assume without loss of generality that I evolves on a finite set.

The input-state dynamics are assumed to be defined by a controlled Markov chain, with controlled transition matrix P. For any randomized stationary policy, the following holds for each  $x, x' \in X$ ,  $u \in U$ , and  $k \ge 0$ :

$$P\{X_{k+1} = x' \mid X_k = x, \ U_k = u\} = P_u(x, x')$$
 (8)

The dynamic programming equation  $Q^* = \mathcal{T}Q^*$  (equivalently (2)) may be expressed, for  $x \in X$ ,  $u \in U$ , by

$$Q^{\star}(x,u) = c(x,u) + \gamma \sum_{x' \in \mathsf{X}} P_u(x,x') \ \underline{Q}^{\star}(x') \tag{9}$$

## B. What is stochastic approximation?

A fuller answer may be found in any of the standard monographs, such as [12] (see also [37] for a crash course).

The goal of SA is to solve the root finding problem  $\bar{f}(\theta^*) = 0$ , where the function is defined in terms of an expectation,  $\bar{f}(\theta) = \mathsf{E}[f(\theta,\Phi)]$  for  $\theta \in \mathbb{R}^d$  and with  $\Phi$  a random vector. The general SA algorithm is expressed in two forms:

$$\theta_{n+1} = \theta_n + \alpha_{n+1} f(\theta_n, \Phi_{n+1})$$
 (10a)

$$= \theta_n + \alpha [\bar{f}(\theta_n) + \Delta_{n+1}], \quad n \ge 0, \tag{10b}$$

where  $\Delta_{n+1} := f(\theta_n, \Phi_{n+1}) - \bar{f}(\theta_n)$ . It is assumed that the sequence  $\{\Phi_n\}$  converges in distribution to  $\Phi$ .

The algorithm is motivated by ordinary differential equation (ODE) theory, and this theory plays a large part in establishing convergence of (10a) along with convergence rates. These results are obtained by comparing solutions (10a) to solutions of the *mean flow*,

$$\frac{d}{dt}\vartheta_t = \bar{f}(\vartheta_t). \tag{11}$$

In particular,  $\theta^*$  is a stationary point of this ODE.

**Averaging** A large step-size  $\{\alpha_{n+1}\}$  in (10a) is desirable for quick transient response, but this typically leads to high variance. There is no conflict if the "noisy" parameter estimates are averaged. The averaging technique of Polyak and Ruppert defines

$$\theta_n^{\text{PR}} = \frac{1}{n} \sum_{k=1}^n \theta_k \,, \qquad n \ge 1. \tag{12}$$

Thm. II.1 illustrates the value of this approach.

**Basic SA assumptions** The following are imposed in this section, and in some others that follow.

It is assumed that the step-size sequence  $\{\alpha_n : n \geq 1\}$  is deterministic, satisfies  $0 < \alpha_n \leq 1$ ,

$$\sum_{n=1}^{\infty} \alpha_n = \infty \quad and \quad \sum_{n=1}^{\infty} \alpha_n^2 < \infty \tag{13}$$

These conditions hold for  $\alpha_n = gn^\rho$  with  $\frac{1}{2} < \rho \le 1$  and g > 0; see [27] for theory justifying larger step-sizes obtained using  $0 < \rho \le \frac{1}{2}$ . We sometimes require two time-scale algorithms in which there is a second step-size sequence  $\{\beta_n : n \ge 1\}$  that is relatively large:

$$\lim_{n \to \infty} \frac{\alpha_n}{\beta_n} = 0 \tag{14}$$

**Parameter dependent noise** A construction of the process  $\Phi$  appearing in the SA recursion reveals one complication. To make the construction entirely explicit, consider a state space realization of the MDP,  $X_{k+1} = F(X_k, U_k, D_k)$  in which D evolves on a finite set, and F is a function taking values in the finite set X. In view of (7) we are in the setting of parameter-dependent noise: instead of one Markov chain, one considers a family  $\{\Phi^{\theta} : \theta \in \mathbb{R}^d\}$ .

In the present setting, the Markov chain  $\Phi^{\theta}$  is defined by freezing the parameter in (7) to define

$$U_k^{\theta} = \phi(X_k^{\theta}, \theta, I_k) \quad X_{k+1}^{\theta} = F(X_k^{\theta}, U_k^{\theta}, D_k), \qquad k \ge 0$$

Note that  $\boldsymbol{X}^{\theta}$  is itself a Markov chain on X, whose transition matrix is denoted  $P_{\theta}$ . A simple choice is then  $\Phi_{k}^{\theta} = (X_{k}, I_{k}, D_{k})$ . Letting  $\xi = (x; \iota; \delta), \, \xi' = (x'; \iota'; \delta')$  denote two state values, the transition matrix is denoted  $\mathcal{P}_{\theta}$  and has the simple realization,

$$\mathcal{P}_{\theta}(\xi, \xi') = P_{\theta}(x, x') \mu_I(\iota') \mu_D(\delta') \tag{15}$$

where  $\mu_I$ ,  $\mu_D$  are the respective pmfs for  $I_k$ ,  $D_k$  (independent of k). It is assumed that  $\mathcal{P}_{\theta}$  admits a unique invariant pmf  $\varpi_{\theta}$  for each  $\theta$ , from which we obtain an expression for the vector field in the mean flow,

$$\bar{f}(\theta) = \mathsf{E}[f(\theta, \Phi^{\theta})], \qquad \Phi^{\theta} \sim \varpi_{\theta}$$
 (16)

That is,  $\Phi^{\theta}$  is distributed according to  $\varpi_{\theta}$  in the expectation.

Theory for convergence of SA with parameter-dependent noise began in the seminal paper [36], and the theory is nearly as mature as in the classical setting with exogenous noise [26]. The following assumptions for the general SA recursion (10) are much stronger than those imposed in this prior work:

## **Assumptions for convergence:**

**SA1** The function f is globally Lipschitz continuous in its first variable, with subgradients satisfying  $\sup_{\theta} |\partial_{\theta_i} f_j(\theta, \xi)| < \infty$  for each i, j and  $\xi$ .

**SA2** For each  $\theta$ , the time-homogeneous Markov chain  $\Phi^{\theta}$  evolves on a finite set. Moreover,

(i) A uniform minorization condition holds: for some  $N \geq 1$ , a constant  $\delta_{\Phi} > 0$ , and a state  $\xi^{\bullet}$ ,

$$\sum_{k=1}^{N} \mathcal{P}_{\theta}^{k}(\xi, \xi^{\bullet}) \ge \delta_{\Phi}, \text{ for each } \theta, \xi.$$
 (17)

It is assumed moreover that  $\Phi^{\theta}$  is aperiodic.

Consequently, there is a unique invariant pmf  $\varpi_{\theta}$ .

(ii) The transition matrix  $\mathcal{P}_{\theta}$  is continuously differentiable in  $\theta$ , with with vanishing gradient for large  $\theta$ : for each  $1 \le i \le d$ ,

$$\sup_{\xi,\xi',\theta} \left| \partial_{\theta_i} \mathcal{P}_{\theta}(\xi,\xi') \right| \|\theta\| < \infty \tag{18}$$

**SA3** The mean flow (11) is globally asymptotically stable, with unique equilibrium  $\theta^*$ .

**Assumptions for convergence rates:** The following is used to obtain useful bounds on the rate of convergence, which requires the existence of a linearization (at least in a neighborhood of  $\theta^*$ ). Denote

$$A(\theta) = \partial_{\theta} \bar{f}(\theta) \tag{19}$$

**SA4** The derivative (19) is a bounded and continuous function of  $\theta$ , and  $A^* := A(\theta^*)$  is a Hurwitz matrix (its eigenvalues lie in the strict left hand plane).

Assumptions (SA1)–(SA3) imply that  $\overline{f}$  is globally Lipschitz continuous. Hence the only part of (SA4) that goes beyond the previous assumptions is the Hurwitz condition.

These assumptions combined with theory in [36], [26] imply convergence of  $\{\theta_n\}$  to  $\theta^*$  almost surely from each initial condition, provided one more property is established:

The parameter sequence  $\{\theta_n : n \ge 0\}$  is bounded with probability one from each initial condition. (20)

Lyapunov techniques provide a means of establishing (20):

(v4) For a globally Lipschitz continuous and  $C^1$  function  $V: \mathbb{R}^d \to [1, \infty)$ , and a constant  $\delta_v > 0$ ,

$$\frac{d}{dt}V(\vartheta_t) \le -\delta_v V(\vartheta_t), \quad \text{when } \|\vartheta_t\| \ge \delta_v^{-1}.$$
 (21)

The designation "v4" comes from an analogous bound appearing in stability theory of Markov chains [39].

An alternative that is often easily verified for RL algorithms is the so-called Borkar-Meyn theorem of [13], [12].

**ODE**@ $\infty$  The time-homogeneous ODE  $\frac{d}{dt}x = \bar{f}_{\infty}(x)$  with vector field,

$$\bar{f}_{\infty}(\theta) := \lim_{r \to \infty} r^{-1} \bar{f}(r\theta).$$
 (22)

We always have  $\bar{f}_{\infty}(0) = 0$ , which means that the origin is an equilibrium for the ODE@ $\infty$ . It is also radially homogeneous,  $\bar{f}_{\infty}(r\theta) = r\bar{f}_{\infty}(\theta)$  for any  $\theta \in \mathbb{R}^d$  and r > 0. Based on these properties it is known that local asymptotic stability of the origin implies global exponential asymptotic stability [13].

Stability of the ODE@ $\infty$  is equivalent to (v4) whenever the limit (22) exits for each  $\theta$ .

It is shown in [13] that (20) holds provided the ODE@ $\infty$  is locally asymptotically stable, and  $\{\Delta_n\}$  appearing in (10b) is a martingale difference sequence. This statistical assumption does *not* hold in many applications of reinforcement learning. Extensions of [13] are given in [9], [41], [10].

The article [10] and its followup [27] require minimal assumptions on the Markov chain (there is no need for a finite state space). While these papers consider exogenous noise, the proof extends to the setting of this paper. See Appendix A for explanation.

**Theorem II.1.** Suppose that (SA1) and (SA2) hold for the SA recursion (10a), and in addition that the origin is locally asymptotically stable for the  $ODE@\infty$ , or that (v4) holds. Then,

(i) The bound (20) holds in a strong sense: there is a fixed constant  $B_{\Theta}$  such that for each initial condition  $(\theta_0, \Phi_0)$ ,

$$\lim_{n \to \infty} \sup \|\theta_n\| \le B_{\Theta} \quad a.s.. \tag{23}$$

(ii) If in addition (SA3) holds then  $\lim_{n\to\infty} \theta_n = \theta^*$  almost surely from each initial condition.

Based on theory in [10] we can expect to also establish mean-square convergence rates—part (iii) that follows is stated as a conjecture at this stage. A functional Central Limit Theorem is easily obtained under the assumptions of Thm. II.1 and also (SA4) [7], [12], implying that the limits in (24) will hold under an  $L_p$  bound on  $\{\tilde{\theta}_n/\sqrt{\alpha_n}:n\geq 1\}$  for some p>2, where the tilde denotes error:  $\tilde{\theta}_n=\theta_n-\theta^*$ ,  $\tilde{\theta}_n^{\rm PR}=\theta_n^{\rm PR}-\theta^*$ . An  $L_4$  bound is established in [10] for exogenous noise.

(iii) Suppose that (SA1)–(SA4) hold, and that  $\alpha_n = gn^\rho$ ,  $n \ge 1$ , with  $\frac{1}{2} < \rho < 1$  and g > 0. We then have convergence in mean square, and the following limits exist and are finite:

$$\lim_{n \to \infty} \frac{1}{\alpha_n} \mathsf{E}[\tilde{\theta}_n \tilde{\theta}_n^{\mathsf{T}}] = \Sigma_{\Theta} \tag{24a}$$

$$\lim_{n \to \infty} n \mathsf{E}[\tilde{\theta}_n^{\mathsf{PR}} \{\tilde{\theta}_n^{\mathsf{PR}}\}^{\mathsf{T}}] = \Sigma_{\Theta}^{\mathsf{PR}}$$
 (24b)

The covariance matrix  $\Sigma_{\Theta}^{PR}$  is minimal in a matricial sense, made precise in [43], [40]. It has the explicit form  $\Sigma_{\Theta}^{PR} = G \Sigma_{\Delta}^* G^{\mathsf{T}}$  in which  $G = -(A^*)^{-1}$ , and

$$\Sigma_{\Delta}^* = \sum_{k=-\infty}^{\infty} \mathsf{E}[\Delta_k^* \{\Delta_k^*\}^{\mathsf{T}}]$$
 (25)

where  $\{\Delta_k^*:=f(\theta^*,\Phi_k^{\theta^*}):k\in\mathbb{Z}\}$ , with  $\Phi^{\theta^*}$  a stationary version of the Markov chain on the two-sided time interval. An alternative representation for  $\Sigma_\Delta^*$  is contained in Appendix A.

In practice we rarely make use of these formulae: the covariance matrix  $\Sigma_{\Theta}^{PR}$  can be estimated using the batch means method, which requires performing many relatively short runs with distinct initial conditions [2].

**A criterion for stationary points** The existence of a suitable Lyapunov function implies the existence of a stationary point.

**Proposition II.2** (Lyapunov Criterion for Existence of a Stationary Point). For an ODE (11) with globally Lipschitz continuous vector field, suppose there is a function  $V: \mathbb{R}^d \to \mathbb{R}_+$  with locally Lipschitz continuous gradient, satisfying for some  $b^{11.2}$ ,

$$\nabla V(\theta)^{\mathsf{T}} \bar{f}(\theta) \leq -1$$
, whenever  $\|\theta\| \geq b^{\mathsf{II}.2}$ .

Suppose moreover that V is convex and coercive. Then there exists a solution to  $\bar{f}(\theta^*) = 0$ .

*Proof.* Let  $L_{\delta}(\theta) = \theta + \delta \bar{f}(\theta)$  for  $\theta \in \mathbb{R}^d$ , with  $\delta > 0$  to be chosen. For  $\delta > 0$  sufficiently small we construct a convex and compact set  $S_{\delta}$  for which  $L_{\delta}(\theta) \in S_{\delta}$  for each  $\theta \in S_{\delta}$ .

It follows from Brouwer's fixed-point theorem that there is a solution to  $L_{\delta}(\theta^*) = \theta^*$ . This is equivalent to the desired conclusion  $\bar{f}(\theta^*) = 0$ .

Denote  $b_{\delta} = \sup\{V(L_{\delta}(\theta)) : \|\theta\| \leq b^{\text{II}2}\}$ , and  $S_{\delta} = \{\theta : V(\theta) \leq b_{\delta}\}$ ; a convex and compact set subject to the assumptions on V.

We next show that  $S_{\delta}$  is invariant under  $L_{\delta}$  if  $\delta$  is small. We consider two cases, based on whether or not  $\theta$  lies in the set  $S = \{\theta : \|\theta\| \le b^{\text{II}.2}\}$ 

- **1.** If  $\theta \in S_{\delta} \cap S$ , then  $L_{\delta}(\theta) \in S_{\delta}$  by construction of  $S_{\delta}$ .
- **2.** If  $\theta \in S_{\delta} \setminus S$  then we apply convexity combined with the drift condition: denoting  $\theta^+ = L_{\delta}(\theta)$ ,

$$V(\theta) \geq V(\theta^+) + \nabla V(\theta^+)^\intercal (\theta - \theta^+) = V(\theta^+) - \delta \nabla V(\theta^+)^\intercal \overline{f}(\theta)$$

Since the gradient is locally Lipschitz continuous and  $\bar{f}$  is globally Lipschitz continuous, there is  $b_v$  satisfying

$$V(\theta) \ge V(\theta^+) - \delta \nabla V(\theta)^{\mathsf{T}} \bar{f}(\theta) - b_v \delta^2, \qquad \theta \in S_\delta \setminus S$$

The value of  $b_v$  can be chosen independent of  $\delta \in (0,1]$ .

Under the assumed drift condition this gives  $V(\theta^+) \le V(\theta) - \delta + b_v \delta^2$ . Choosing  $\delta = 1/b_v$  gives  $V(\theta^+) \le V(\theta) \le b_\delta$ , in which the second inequality holds because  $\theta \in S_\delta \setminus S$ . Hence  $L_\delta(\theta) = \theta^+ \in S_\delta$  as desired.

When stability of the mean flow cannot be established, stability can typically be assured using a matrix gain algorithm.

**Zap stochastic approximation.** This is a two time-scale algorithm introduced in [20]. For initialization  $\theta_0 \in \mathbb{R}^d$ , and  $\widehat{A}_0 \in \mathbb{R}^{d \times d}$ , obtain the sequence of estimates  $\{\theta_n : n \geq 0\}$  recursively:

$$\theta_{n+1} = \theta_n - \alpha_{n+1} \widehat{A}_{n+1}^{-1} f(\theta_n, \Phi_{n+1})$$
 (26a)

$$\widehat{A}_{n+1} = \widehat{A}_n + \beta_{n+1} [A_{n+1} - \widehat{A}_n],$$
 (26b)

with  $A_{n+1} := \partial_{\theta} f_{n+1}(\theta_n)$ .

The two gain sequences  $\{\alpha_n\}$  and  $\{\beta_n\}$  satisfy (14). This ensures that the ODE approximation for the parameter estimates is the Newton-Raphson flow  $\frac{d}{dt}\bar{f}(\vartheta) = -\bar{f}(\vartheta)$ . Hence stability is assured if  $\|\bar{f}\|$  is coercive [14].

The recursion (26b) requires modification for the applications considered here, in which the transition law for the Markov chain depends on the parameter estimate. See discussion surrounding eq. (51) in Section IV-D.

#### C. Compatible assumptions for Q-learning

The basic Q-learning algorithm (6a) is an instance of stochastic approximation, for which we can apply general theory subject to assumptions on the input used for training (recall (7)). Two settings are considered:

**Oblivious training** This means that (7) simplifies to

$$U_k = \phi(X_k, I_k), \qquad k \ge 0, \tag{27}$$

in which it is always assumed that  $\{I_k\}$  is i.i.d..

It follows that the pair process  $\{(X_k, U_k) : k \geq 0\}$  is a time homogeneous Markov chain. It is assumed to be unichain (i.e., the invariant pmf  $\pi$  is unique). In the expression  $f_{n+1}(\theta_n) = f(\theta_n, \Phi_{n+1})$  we take  $\{\Phi_k = (X_k; X_{k+1}; U_k) : \Phi_n\}$ 

 $k \ge 0$ }, which is also a time homogeneous Markov chain, for which its invariant pmf is also unique and easily expressed in terms of  $\pi$  and the controlled transition matrix.

If the function class is linear  $\{Q^{\theta} = \theta^{\mathsf{T}} \psi : \theta \in \mathbb{R}^d\}$ , then the autocorrelation matrix is assumed full rank

$$R_0 = \mathsf{E}_{\pi}[\psi(X_n, U_n)\psi(X_n, U_n)^{\mathsf{T}}] \tag{28}$$

where the expectation is taken in steady-state

**Optimistic training** In this non-oblivious approach the input sequence depends on the parameter sequence, and is designed to approximate the  $Q^{\theta}$ -greedy policy  $\Phi^{\theta}$  defined in (3).

There are only a finite number of deterministic stationary policies, so  $\phi^{\theta}$  is necessarily discontinuous in  $\theta$ . The region on which continuity holds is denoted

$$\mathcal{C}^{\Theta} = \begin{cases} \theta \in \mathbb{R}^d : \text{there is } \varepsilon > 0 \text{ s.t. } \varphi^{\theta}(x) = \varphi^{\theta'}(x) \\ \text{for all } x \text{ when } \|\theta - \theta'\| \le \varepsilon \end{cases}$$
 (29)

The training policy is taken of the form,

$$U_k = (1 - B_k)\mathcal{U}_k + B_k\mathcal{W}_k \tag{30}$$

in which  $\{B_k\}$  is an i.i.d. Bernoulli sequence with  $P\{B_k = 1\} = \varepsilon$ , and  $\{W_k\}$  is an i.i.d. sequence taking values in U and independent of  $\{B_k\}$ . The U-valued random variable  $\mathcal{U}_k$  depends on the parameter  $\theta_k$ , and is independent of  $(B_k; \mathcal{W}_k)$  for each k.

The sequences  $\{U_k, \mathcal{U}_k : k \geq 0\}$  are defined by randomized stationary policies  $\{\widetilde{\Phi}^{\theta}, \widetilde{\Phi}^{\theta}_0 : \theta \in \mathbb{R}^d\}$ . Both  $\widetilde{\Phi}^{\theta}(\cdot \mid x)$  and  $\widetilde{\Phi}^{\theta}_0(\cdot \mid x)$  are pmfs on U for each x and  $\theta$ . Based on the assumptions imposed after (30), we have

$$P\{U_{k} = u \mid \mathcal{F}_{k}^{-}; X_{k} = x\} = \widetilde{\Phi}^{\theta_{k}}(u \mid x)$$
$$= (1 - \varepsilon)\widetilde{\Phi}_{0}^{\theta_{k}}(u \mid x) + \varepsilon \nu_{w}(u)$$
(31)

with  $v_{\mathcal{W}}$  the common pmf for  $\{\mathcal{W}_k\}$ , and  $\mathcal{F}_k^- = \sigma\{X_i, U_i : i < k; B_i, \mathcal{W}_i : i \le k\}$  (a partial history of observations up to iteration k).

Special cases are described in the following.

**1.**  $\varepsilon$ **-greedy.** The choice  $\mathcal{U}_k = \phi^{\theta_k}(X_k)$ , so that

$$\widetilde{\Phi}_0^{\theta}(u \mid x) = \mathbb{1}\{u = \Phi^{\theta}(x)\}\tag{32}$$

The mean flow has many attractive properties (see Prop. A.6 in the Appendix). However, because  $\{\Phi^{\theta}: \theta \in \mathbb{R}^d\}$  is a piecewise constant function of  $\theta$ , it follows that the vector field  $\bar{f}$  is not continuous in  $\theta$  as required in Thm. II.1.

**2. Gibbs approximation** For fixed constant  $\kappa > 0$  define

$$\widetilde{\Phi}_0^{\theta}(u \mid x) = \frac{1}{\mathcal{Z}_{\kappa}^{\theta}(x)} \exp(-\kappa Q^{\theta}(x, u))$$
 (33)

in which  $\mathcal{Z}_{\kappa}^{\theta}(x)$  is normalization. This is indeed an approximation of (32): for  $\theta \in \mathcal{C}^{\Theta}$ ,

$$\lim_{r \to \infty} \frac{1}{\mathcal{Z}_{\kappa}^{r\theta}(x)} \exp(-\kappa Q^{r\theta}(x, u)) = \mathbb{1}\{u = \phi^{\theta}(x)\}$$
 (34)

The limit (34) has two important implications. First is that the vector field  $\bar{f}_{\infty}$  for the ODE@ $\infty$  is unchanged whether we consider (32) or its smooth approximation (33). Second is that discontinuity of  $\bar{f}_{\infty}$  implies that  $\bar{f}$  is not globally Lipschitz continuous, which violates an assumption of Thm. II.1.

**3. Tamed Gibbs approximation** This is a modification of (33) in which  $\kappa$  depends on  $\theta$ :

$$\widetilde{\Phi}_0^{\theta}(u \mid x) = \frac{1}{\mathcal{Z}_{\kappa}^{\theta}(x)} \exp(-\kappa_{\theta} Q^{\theta}(x, u))$$
 (35)

For analysis the following structure is helpful: choose a large constant  $\kappa_0 > 0$ , and assume that

$$\kappa_{\theta} \begin{cases} = \frac{1}{\|\theta\|} \kappa_{0} & \|\theta\| \ge 1 \\ \ge \frac{1}{2} \kappa_{0} & else \end{cases}$$
 (36)

This will be called the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy when it is necessary to make the policy parameters explicit.

The equality in (36) ensures the following holds for all x, u:

$$\widetilde{\Phi}^{r\theta}(u \mid x) = \widetilde{\Phi}^{\theta}(u \mid x) \text{ for all } r \ge 1 \text{ and } \|\theta\| \ge 1.$$
 (37)

The Q-learning algorithm (6) can be cast as stochastic approximation when the input is defined using any of the training policies described above, in which we take  $\Phi_{n+1} = (X_n, X_{n+1}, U_n)$  since these three variables appear in (6).

It is assumed in Thm. II.1 that  $\Phi$  is *exogenous*—its transition matrix does not depend on the parameter sequence. Fortunately, there is now well developed theory that allows for parameter-dependent dynamics for  $\Phi$  in the SA recursion (10a)—see the recent paper [56] for history and recent results. In particular, theory of convergence and asymptotic statistics is now mature.

The question is then, how can we apply SA theory to make statements about convergence and convergence rates?

#### III. TROUBLE WITH TABULAR

We begin with a useful representation for the mean flow vector field  $\bar{f}$  for Q-learning with linear function approximation (5), for arbitrary basis. The projected Bellman equation (4) is the root finding problem,  $\bar{f}(\theta^*) = 0$ .

To avoid long equations we adopt the shorthand notation,

$$f_{n+1}(\theta_n) = f(\theta_n, \Phi_{n+1}) c_n = c(X_n, U_n), \quad \psi_{(n)} = \psi(X_n, U_n).$$
 (38)

The eligibility vector in (6a) is then  $\zeta_n = \psi_{(n)}$ .

If the parameter  $\theta$  is frozen, so that  $U_k \sim \widetilde{\Phi}^{\theta}(\cdot \mid X_k)$  for each k, then the controlled state process  $\boldsymbol{X}^{\theta}$  is a time homogeneous Markov chain with transition matrix,

$$P_{\theta}(x, x') := \sum_{u} \widetilde{\Phi}^{\theta}(u \mid x) P_{u}(x, x'), \qquad x, x' \in X. \quad (39a)$$

The pair process  $\{(X_k, U_k) : k \ge 0\}$  is also Markovian, with transition matrix denoted

$$T_{\theta}(z, z') := P_{u}(x, x')\widetilde{\Phi}^{\theta}(u' \mid x'), \qquad (39b)$$

for each  $z=(x,u)\,,\;z'=(x',u')\in\mathsf{X}\times\mathsf{U}.$  It is assumed that  $T_{\theta}$  has a unique invariant pmf  $\pi_{\theta}.$ 

Of course, the parameter  $\theta$  is never frozen in any algorithm. The transition matrices  $P_{\theta}$  and  $T_{\theta}$  are introduced for analysis.

Q-learning in the form (6a) is an instance of stochastic approximation, with mean flow vector field,

$$\bar{f}(\theta) = \mathsf{E}_{\pi_{\theta}}[\psi_{(n)}\mathcal{B}(X_n, U_n; \theta)], \qquad (40a)$$

 $\mathcal{B}(x, u; \theta) = c(x, u) - Q^{\theta}(x, u)$ 

$$+\gamma \sum_{x'} P_u(x, x') \underline{Q}^{\theta}(x')$$
 (40b)

An alternative formula is valuable for analysis.

Lemma III.1. The vector field (40a) may be expressed

$$\bar{f}(\theta) = A(\theta)\theta - b(\theta) 
with A(\theta) = -\mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n)} - \gamma \psi_{(n+1)}^{\theta} \}^{\mathsf{T}} \right] 
b(\theta) = -\mathsf{E}_{\pi_{\theta}} [\psi_{(n)} c_n]$$
(41)

and 
$$\psi_{(n+1)}^{\theta} = \psi(X_{n+1}, u)$$
 with  $u = \Phi^{\theta}(X_{n+1})$ .

The vector field  $\bar{f}$  is globally Lipschitz continuous when using the training policy (31) with tamed Gibbs policy (35), for any value of  $\varepsilon \in [0,1]$  and  $\kappa > 0$ .

The representation (41) follows directly from (40b). Lipschitz continuity follows Lemma A.1 combined with Prop. A.2 (each postponed to the Appendix).

Note that  $\varepsilon=1$  corresponds to an oblivious policy, so the lemma provides a large collection of policies for which  $\bar{f}$  fits the standard SA theory. The tamed Gibbs policy is the only choice among the optimistic training rules for which  $\bar{f}$  satisfies the smoothness conditions required in Thm. II.1.

In the remainder of this section we restrict to oblivious training. The main results of this paper in Section IV concern optimistic training.

A. Tabular Q-learning, the good and the bad

In the tabular setting we take  $d = |X| \times |U|$  in (5), and

$$\psi_i(x,u) = \mathbb{1}\{(x,u) = (x^i,u^i)\}, \qquad x \in X, \ u \in U \quad (42)$$

where  $\{(x^i, u^i) : 1 \le i \le d\}$  is any ordering of state-action pairs. Hence  $Q^{\theta_n}(x^i, u^i) = \theta_n(i)$  for each n, i.

It is typical to use a diagonal matrix gain,

$$\theta_{n+1} = \theta_n + \alpha_{n+1} G_n \mathcal{D}_{n+1} \zeta_n \tag{43}$$

in which  $G_n^{-1}(i,i)$  indicates the number of times the pair  $(x^i, u^i)$  is visited up to time n (set to unity when this is zero).

The mean flow (11) associated with (43) is

$$\frac{d}{dt}\vartheta_t = A(\vartheta_t)\vartheta_t - b \tag{44}$$

with b the d-dimensional vector with entries  $b_i = -c(x^i, u^i)$ , and the matrix-valued function A is piecewise constant.

**The good news:** The statistical properties of the algorithm are attractive because  $\{\Delta_{n+1}\}$  appearing in (10b) is a martingale difference sequence in the tabular setting.

The best news is stability: we have  $A(\theta) = -[I - \gamma M(\theta)]$ , in which  $M_{i,j}(\theta) = P_{u^i}(x^i, x^j) \mathbb{1}\{u^j = \Phi^{\theta}(x^j)\}$ . The induced operator norm of  $M(\theta)$  in  $\ell_{\infty}$  is no greater than one, meaning  $\max_i |\sum_j M_{i,j}(\theta) v_j| \leq ||v||_{\infty} := \max_i |v_i|$  for any vector v

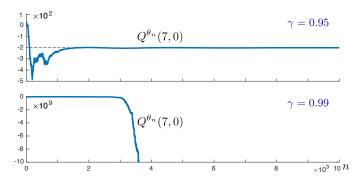


Fig. 1. Evolution of the Q-function approximations for two values of discount factor, and using an  $\varepsilon$ -greedy policy with common value of  $\varepsilon = 0.5$ .

and any  $\theta$ . It follows that the  $\ell_{\infty}$  norm serves as a Lyapunov function: Letting  $\tilde{\vartheta}_t = \vartheta_t - \theta^*$  and  $V(\theta) = \|\theta\|_{\infty}$ ,

$$\frac{d}{dt}V(\vartheta_t) \le -(1-\gamma)V(\vartheta_t)$$

This is how convergence is established for tabular Q-learning. **The bad:** The matrix A(q) has an eigenvalue at  $-(1-\gamma)$  for each  $\theta$ , which is a reason for slow convergence when the discount factor is close to unity. One consequence is that the asymptotic covariance  $\Sigma_{\Theta}$  appearing in (24a) is not finite if  $\gamma > 1/2$  and the step-size sequence is  $\alpha_n = 1/n$  (see [20] and the sample complexity analysis that followed in [52]).

#### B. Change your goals

A reader with experience in SA would counter that  $\alpha_n=1/n$  is a poor choice of step-size. Use instead  $\alpha_n=1/n^\rho$ , with  $\rho\in(\frac{1}{2},1)$ , and then average using (12) to obtain  $\{\theta_n^{\rm PR}\}$ . It is found that averaging fails for this example for large discount factors, even though it is known that these estimates achieve the optimal asymptotic covariance [37], [19], [18].

The observed numerical instability is a consequence of the eigenvalue at  $-(1-\gamma)$  for  $A^*:=A(\theta^*)$  (recall (19)). The eigenvalue can be moved through a change in objective. For example, construct an algorithm that estimates the *relative Q-function*,

$$H^{\star}(x,u) = Q^{\star}(x,u) - \delta\langle \mathbf{v}, Q^{\star} \rangle \tag{45}$$

where  $\nu$  is a fixed pmf on X × U and  $\delta$  > 0. Subtracting a constant doesn't change the minimizer over u, and has enormous benefits.

The function  $H^*$  satisfies a DP equation which motivates relative Q-learning. It is shown in [21] that the eigenvalues of  $A^*$  remain bounded away from the imaginary axis uniformly for all  $0 \le \gamma \le 1$ , resulting in much faster convergence. See [37] for generalizations.

## IV. STABILITY WITH OPTIMISM

The theory surveyed in the preceding section imposed oblivious training. In the case of Watkins' Q-learning this assumption was imposed in part for historical reasons, though we will see that the analysis is somewhat more complex when we consider parameter dependent policies. The technical challenges for Zap Q-learning are far more interesting because

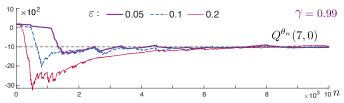


Fig. 2. Evolution of the Q-function approximations when using an  $\varepsilon$ -greedy policy. Convergence holds when  $\varepsilon > 0$  is sufficiently small.

the definition of the linearization  $A(\theta)$  is not obvious. See the conclusions for further discussion.

We begin with a motivating example.

#### A. Baird's star example

We refer the reader to the source [4]—the final page contains a full description of the model considered in the experiments surveyed here. See [45] for a fuller discussion.

There are seven states  $X = \{1, \ldots, 7\}$  and two actions  $U = \{0, 1\}$ , in which  $X_{k+1} = 7$  with probability one whenever  $U_k = 0$ . In [4] it is assumed the cost is identically zero. We take c(x, u) = 0 if  $x \le 6$  and c(7, u) = -10 (independent of u). The Q-function is linearly parameterized with dimension d = 14. With a well-motivated oblivious policy it was shown that the parameter estimates diverge when the discount factor is sufficiently large.

Fig. 1 shows trajectories from the Q-learning algorithm (6a) with an  $\varepsilon$ -greedy policy using  $\varepsilon=0.5$ . The ideal behavior is that  $Q^{\theta_n}(x,u) \to Q^{\star}(x,u) = -10/(1-\gamma)$  as  $n \to \infty$  when (x,u)=(7,0). The figure shows convergence when  $\gamma=0.95$ , but the parameters are divergent with discount factor  $\gamma=0.99$ .

With the larger discount factor we obtain stability when using a smaller value of  $\varepsilon > 0$ . Fig. 2 shows typical results for three small values. The dashed line indicates  $Q^*(7,0)$ .

The step-size sequence was taken to be  $\alpha_n = \min(\bar{\alpha}, g/n^{\rho})$  using  $g = 1/(1-\gamma)$ ,  $\rho = 0.85$ , and  $\bar{\alpha} = 0.1$  in each run. The Matlab code is available on arXiv—see the Appendix of [38].

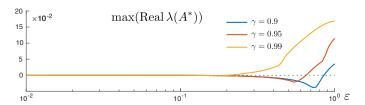


Fig. 3. The maximum eigenvalue of  $A^*$  as a function of  $\varepsilon$ . The matrix is Hurwitz for sufficiently small  $\varepsilon>0$ , but some eigenvalues approach zero with vanishing  $\varepsilon$ .

See Section III-B for an explanation for slow convergence with a large discount factor, and [21] for explanation of the choice  $g=1/(1-\gamma)$  based on consideration of the linearization matrix  $A^*$ . Fig. 3 shows a plot of the maximum real part of  $A^*$  as a function of  $\varepsilon>0$ , estimated via Monte-Carlo. For larger values of  $\varepsilon>0$  we see that  $A^*$  is not Hurwitz for the three choices of discount factor. There is also trouble for very small  $\varepsilon>0$ : The discussion following Thm. II.1 suggests that the asymptotic covariance will be very large

when  $\max(\operatorname{Real} \lambda(A^*))$  is close to zero, but the covariance  $\Sigma_{\Delta}^*$  must also be considered to make any conclusions.

Applications to change detection Similar experiments were conducted in [16] for application to quickest change detection. Much like in Baird's example it was found that  $\varepsilon$ -greedy training was successful, but only for extremely small values of  $\varepsilon>0$ . Zap Q-learning was far more reliable over a large range of  $\varepsilon\in(0,1)$ : testing revealed that the final parameter  $\theta^{\bullet}$  defining the policy was nearly optimal. However, the matrix  $A(\theta^{\bullet})$  was not Hurwitz, which suggests that the standard algorithm (6) is unable to recover  $\theta^{\bullet}$ .

These findings illustrate that significant understanding of RL theory is essential for practical success in many applications.

## B. Sufficient optimism

The main result of this paper shows how exploration using a policy of the form (30) encourages stability of the Q-learning algorithm (6) with linear function approximation (5). Analysis of (6) requires the family of autocorrelation matrices,

$$R^{\Theta}(\theta) = \mathsf{E}_{\pi_{\theta}} \left[ \psi(X_n, \phi^{\theta}(X_n)) \psi(X_n, \phi^{\theta}(X_n))^{\mathsf{T}} \right] \tag{46a}$$

$$R^{\mathcal{W}}(\theta) = \mathsf{E}_{\pi_{\theta}} \left[ \psi(X_n, \mathcal{W}_n) \psi(X_n, \mathcal{W}_n)^{\mathsf{T}} \right] \tag{46b}$$

$$R(\theta) = \mathsf{E}_{\pi_{\theta}}[\psi_{(n)}\psi_{(n)}^{\mathsf{T}}] = (1 - \varepsilon)R^{\Theta}(\theta) + \varepsilon R^{\mathcal{W}}(\theta) \quad (46c)$$

The expectations are in steady-state, with stationary pmf  $\pi_{\theta}$  induced by the randomized stationary policy with fixed parameter.

A special case is considered in the assumptions, in which we take  $\varepsilon=1$ , and the randomized policy is then denoted  $\widetilde{\Phi}^{\mathcal{W}}$ , giving  $\widetilde{\Phi}^{\mathcal{W}}(\cdot\mid x)=\nu_{\mathcal{W}}(u)$  for all x,u. The (assumed unique) invariant pmf is denoted  $\pi_{\mathcal{W}}$ , and the autocorrelation matrix

$$R^{\mathcal{W}} = \mathsf{E}_{\pi_{\mathcal{W}}} \big[ \psi(X_n, \mathcal{W}_n) \psi(X_n, \mathcal{W}_n)^{\mathsf{T}} \big]$$

$$using \ U_k = \mathcal{W}_k \ for \ all \ k.$$
(46d)

We have  $R^{w} > 0$  in Baird's star example whenever the distribution of  $W_k$  is not degenerate.

The following assumptions are required in the main results of this section:

The randomized policy  $\widetilde{\Phi}^{\mathcal{W}}$  gives rise to an aperiodic and uni-chain Markov chain, with unique invariant pmf  $\pi_{\mathcal{W}}$ , and the autocorrelation matrix  $R^{\mathcal{W}}$  defined in (46d) is positive definite. (47a)

The inverse temperature  $\kappa_{\theta}$  is twice continuously differentiable ( $C^2$ ) in  $\theta$ , and the first and second derivatives of  $\kappa_{\theta}$  are continuous and bounded. (47b)

We also require small  $\varepsilon>0$  in specification of the policies. Denote

$$\varepsilon_{\gamma} := \frac{(1-\gamma)^2}{(1-\gamma)^2 + \gamma^2} \tag{48}$$

**Theorem IV.1.** Consider the Q-learning algorithm (6a) with linear function approximation, and training policy (30) defined using the tamed Gibbs policy (35). Suppose moreover that (47) holds. Then, for any  $\varepsilon \in (0, \varepsilon_{\gamma})$  there is  $\kappa_{\varepsilon, \gamma} > 0$  for which the following hold using the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy, using  $\kappa_0 \geq \kappa_{\varepsilon, \gamma}$ :

- (i) The parameter estimates  $\{\theta_n\}$  are bounded: there is a fixed constant  $B_{\Theta}$ , independent of  $\kappa_0 \geq \kappa_{\varepsilon,\gamma}$ , such that (23) holds with probability one from each initial condition.
- (ii) There exists at least one solution to the projected Bellman equation (4).

See Section IV-C for an extension of (ii) to the  $\varepsilon$ -greedy policy.

To see why (i) is plausible, consider an algorithm approximating (6a), in which the minimum defining  $\underline{Q}^{\theta_n}(X_{n+1})$  is replaced by substitution of the input used for training:

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \widetilde{\mathcal{D}}_{n+1} \zeta_n .$$

$$\widetilde{\mathcal{D}}_{n+1} = c_n - Q^{\theta_n}(X_n, U_n) + \gamma Q^{\theta_n}(X_{n+1}, U_{n+1}^-)$$
(49)

in which  $U_{n+1}^-$  is obtained by sampling from  $\widetilde{\Phi}^{\theta}(\cdot \mid x)$  using  $x = X_{n+1}$  and  $\theta = \theta_n$ . The recursion (49) is a variant of the SARSA algorithm [42], [45].

Stability of the ODE@ $\infty$  is then relatively easy, from which we obtain the following:

**Proposition IV.2.** Consider the recursion (49) with linear function approximation, and training policy (30) defined as the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy (35) with  $\varepsilon \in (0, 1)$  and  $\kappa_0 > 0$ . Suppose moreover that (47) holds. Then, we obtain the conclusions of Thm. IV.1:

- (i) The parameter estimates  $\{\theta_n\}$  are bounded with probability one from each initial condition.
- (ii) There exists at least one solution  $\theta^*$  to  $\bar{f}(\theta^*) = 0$ , with  $\bar{f}$  the mean flow for (49).

The proof of Thm. IV.1 is postponed to the Appendix—we proceed here with the proof of Prop. IV.2.

To begin, suppose that  $U_{n+1}^-$  is replaced by  $U_{n+1}$  in (49). This does not lead to a practical algorithm, since  $\widetilde{\mathcal{D}}_{n+1}$  would then depend on  $\theta_{n+1}$ , but it may be regarded as an approximation since  $\theta_{n+1} \approx \theta_n$ . The approximation leads to a recursion similar to the TD(0) learning algorithm:

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \left[ \psi_{(n)} c_n - \psi_{(n)} \{ \psi_{(n)} - \gamma \psi_{(n+1)} \}^{\mathsf{T}} \theta_n \right]$$

This motivates consideration of the family of autocorrelation matrices  $R_k(\theta) = \mathsf{E}_{\pi_\theta}[\psi_{(n+k)}\psi_{(n)}^\intercal]$  for  $n,k \geq 0$ , so that  $R_0(\theta) = R(\theta)$  in the notation (46c).

The vector field for the mean flow associated with (49) is Lipschitz continuous and has an attractive form in terms of the vector and matrix valued functions,

$$b(\theta) = -\mathsf{E}_{\pi_{\theta}}[\psi_{(n)}c_n], \quad A(\theta) = -R_0(\theta) + \gamma R_{-1}(\theta)$$

**Lemma IV.3.** Under the assumptions of Prop. IV.2 the following hold for (49):

- (i) The vector field for the mean flow is  $\bar{f}(\theta) = A(\theta)\theta b(\theta)$ .
- (ii) The limit defining  $\bar{f}_{\infty}$  in (22) exists and may be expressed  $\bar{f}_{\infty}(\theta) = A_{\infty}(\theta)\theta$  where  $A_{\infty}(\theta) = A(\theta/\|\theta\|)$  for  $\theta \neq 0$ .

*Proof.* Identification of  $\bar{f}$  follows immediately from (49) since  $\theta$  is held fixed in the definition of the mean flow. The representation of the ODE@ $\infty$  follows from structure of the

policy highlighted in (37), which implies the following for all  $r \geq 1$  and  $\theta \in \mathbb{R}^d$  satisfying  $\|\theta\| \geq 1$ :

$$\pi_{r\theta} = \pi_{\theta}$$
,  $A(r\theta) = A(\theta)$ , and  $b(r\theta) = b(\theta)$ 

**Lemma IV.4.** Suppose that (47a) holds. Then, for the recursion (49) there exists  $\delta_{\psi} > 0$ , independent of  $\theta$  such that

$$R_0(\theta) \ge \delta_{\psi} I$$
 for all  $\theta \in \mathbb{R}^d$   
 $\theta^{\mathsf{T}} A(\theta) \theta < -(1-\gamma)\delta_{\psi}$  for all  $\theta \in \mathbb{R}^d$ ,  $\|\theta\| > 1$ .

*Proof.* The proof of the lower bound on  $R_0(\theta)$  is identical to the proof of Lemma A.3 in the Appendix. From Lemma IV.3 (i) we have for  $\theta \in \mathbb{R}^d$  satisfying  $\|\theta\| > 1$ ,

$$\begin{split} \theta^\intercal A(\theta) \theta &= -\theta^\intercal R_0(\theta) \theta + \gamma \theta^\intercal R_{-1}(\theta) \theta \\ &\leq -(1-\gamma) \theta^\intercal R_0(\theta) \theta \leq -(1-\gamma) \delta_\psi \|\theta\|^2 \end{split}$$

*Proof of Prop. IV.2.* Let  $V_1(\theta) = \frac{1}{2} \|\theta\|^2$  and apply Lemmas IV.3 and IV.4 to obtain, whenever  $\|\vartheta_t\| \geq 1$ ,

$$\frac{d}{dt}V_1(\vartheta_t) = \vartheta_t^{\mathsf{T}} \bar{f}(\vartheta_t) = \vartheta_t^{\mathsf{T}} \{A(\vartheta_t)\vartheta_t - b(\vartheta_t)\} 
\leq -\delta_1 \|\vartheta_t\|^2 + \|\vartheta_t\| \|b(\vartheta_t)\|$$

with  $\delta_1 = (1 - \gamma)\delta_{\psi}$ . This gives, with  $\bar{b} = \sup_{\theta} \|b(\theta)\| < \infty$ ,

$$\frac{d}{dt}V_1(\vartheta_t) \le -\frac{1}{2}\delta_1 \|\vartheta_t\|^2, \qquad \|\vartheta_t\| \ge \max(1, 2\bar{b})$$

We then obtain (v4) using  $V(\theta) = \sqrt{V(\theta)} = \|\theta\|$  for  $\|\theta\| \ge \max(1, 2\bar{b})$  (modified in a neighborhood of the origin to impose the  $C^1$  condition):

$$\frac{d}{dt}V(\vartheta_t) \leq -\delta_v V(\vartheta_t), \qquad \|\vartheta_t\| \geq \max(1, 2\bar{b}),$$

with  $\delta_v = \delta_1/4$ . Part (i) then follows from Thm. II.1 (i) and part (ii) from Prop. II.2.

## C. Implications to the $\varepsilon$ -greedy policy

A full analysis of Q-learning using the  $\varepsilon$ -greedy policy for training is beyond the scope of this paper due to discontinuity of the vector field. We find here that Thm. IV.1 admits a partial extension.

We consider here the mean flow (40a), and also the algorithm with matrix gain, whose mean flow vector field is

$$\bar{f}^{\mathrm{zap}}(\theta) = -[A(\theta)]^{-1}\bar{f}(\theta) = -\theta + [A(\theta)]^{-1}b(\theta)\,, \quad \theta \in \mathcal{C}^{\ominus}$$

This defines the dynamics expected when using Zap Q-learning based on (26).

The set  $C^{\Theta}$  in (29) may be expressed as the disjoint union,

$$\mathcal{C}^{\scriptscriptstyle \Theta} = \bigcup_{i} \mathcal{C}^{\scriptscriptstyle \Theta}_i$$

in which each  $\mathcal{C}_i^{\scriptscriptstyle \Theta}$  is an open convex polyhedron, with  $\varphi^{\theta}=\varphi^{\theta'}$  for all  $\theta,\theta'\in\mathcal{C}_i^{\scriptscriptstyle \Theta}$ . Consequently, both  $\bar{f}$  and  $\bar{f}^{\sf zap}$  are constant on each set  $\mathcal{C}_i^{\scriptscriptstyle \Theta}$ .

For each  $\theta \in \mathbb{R}^d$ , denote by  $\Phi^{\theta}$  the set of all randomized  $Q^{\theta}$ -greedy policies: if  $\widetilde{\Phi} \in \Phi^{\theta}$  then

$$\sum_{u} \widetilde{\phi}(u \mid x) Q^{\theta}(x, u) = \underline{Q}^{\theta}(x), \qquad x \in X.$$

If  $\theta \in \mathcal{C}^{\Theta}$  then  $\Phi^{\theta} = \{ \varphi^{\theta} \}$  is a singleton.

**Theorem IV.5.** Suppose that (47a) holds. Then, the following hold for the mean flows associated with the Q-learning algorithm with  $\varepsilon$ -greedy training, provided  $0 < \varepsilon < \varepsilon_{\gamma}$ :

(i) There exists  $\theta^* \in \mathbb{R}^d$  and  $\widetilde{\Phi}^* \in \Phi^{\theta^*}$  such that  $\overline{f}(\theta^*) = 0$ , with  $\overline{f}$  defined in (40a) in which the expectation is taken in steady-state using  $\pi_{\theta^*}$  obtained from the randomized policy,

$$\widetilde{\Phi}^{\theta^*}(u \mid x) = (1 - \varepsilon)\widetilde{\Phi}^*(u \mid x) + \varepsilon \nu_{\mathcal{W}}(u) \tag{50}$$

- (ii) If  $\theta^* \in C^{\Theta}$  then  $\theta^*$  is locally asymptotically stable for the mean flow with vector field  $\bar{f}$ .
- (iii) If  $\theta^* \in C_i^{\Theta}$  for some i, then  $\theta^*$  is locally asymptotically stable for the mean flow with vector field  $\bar{f}^{\mathsf{zap}}$ , with domain of attraction including all of  $C_i^{\Theta}$ .

*Proof.* The proof of (i) is contained in Appendix E.

If  $\bar{f}(\theta^*) = 0$  with  $\theta^* \in \mathcal{C}^{\Theta}$ , it then follows from the definition of the vector field that  $\theta^* = [A(\theta^*)]^{-1}b(\theta^*)$ . Consequently, for  $\theta$  in a neighborhood of  $\theta^*$  contained in  $\mathcal{C}^{\Theta}$ ,

$$\bar{f}(\theta) = A(\theta^*)(\theta - \theta^*)$$

See Prop. A.6 for a proof that  $A(\theta^*)$  is Hurwitz, so that  $\theta^*$  is locally asymptotically stable as claimed in (ii).

We have under the assumptions of (iii),

$$\bar{f}^{\mathsf{zap}}(\theta) = -\theta + \theta^* \,, \quad \theta \in \mathcal{C}_i^{\scriptscriptstyle \Theta}$$

If  $\vartheta_0 \in \mathcal{C}_i^{\Theta}$  it follows that the solution to  $\frac{d}{dt}\vartheta_t = \bar{f}^{\mathsf{zap}}(\vartheta_t)$  is given by  $\vartheta_t = \theta^* + [\vartheta_t - \theta^*]e^{-t}$ . Convexity of  $\mathcal{C}_i^{\Theta}$  ensures that  $\vartheta_t \in \mathcal{C}_i^{\Theta}$  for all t, which completes the proof of (iii).

## D. Extensions from the basic algorithm

Theory for the basic Q-learning algorithm will have implications to other algorithms:

**Stochastic Gradient Descent (SGD)** The mean flow for the GQ learning algorithm of [31] can be expressed  $\frac{d}{dt}\vartheta_t = -\nabla_\theta L\left(\vartheta_t\right)$ , with  $L(\theta) = \bar{f}(\theta)^\intercal Z\bar{f}(\theta)$  and Z positive definite. The theory in the present paper provides sufficient conditions under which the non-singularity condition (L3) of [31] holds, and most important is the new finding that  $\min_\theta L(\theta) = 0$ .

A numerical challenge with GQ learning or gradient descent [4], [3] is that the condition number of the linearization is squared. In particular, for GQ-learning the linearization is expressed  $\frac{d}{dt}\tilde{\vartheta}_t \approx -[\nabla_{\theta}^2 L\left(\theta^*\right)]\tilde{\vartheta}_t$ , and  $\lambda_{\min}(\nabla_{\theta}^2 L\left(\theta^*\right)) = O(|1-\gamma|^2)$ . Hence some of the bad news reviewed in Section III-A is exacerbated using these SGD methods.

**Relative Q-learning** The mean flow vector field for this algorithm is a modification of (41):  $\bar{f}(\theta) = [A(\theta) - Z]\theta - b(\theta)$ , in which Z is a rank one matrix chosen by the user (the specific form follows from (45)). [37, Proposition 9.23] (adapted from [21]) tells us that that the maximum eigenvalue of  $A^*$  remains bounded away from 0 for relative Q-learning in the tabular setting. It is conjectured that Thm. IV.1 can be extended to these algorithms, with  $\varepsilon_{\gamma}$  sufficiently small, but independent of the discount factor  $\gamma \in (0,1)$ . This may require a fresh look at the choice of Z.

**Regularization** Similar to relative Q-learning, in the regularized Q-learning algorithm of [29] the mean flow becomes  $\bar{f}(\theta) = [A(\theta) - Z]\theta - b(\theta)$ . It is possible that the conditions on Z for convergence may be relaxed based on the theory in this paper.

**Double Q-learning** Stability has been established in the tabular setting [23]. The algorithm with linear function approximation has a 2d-dimensional mean flow whose state  $\vartheta_t = [\vartheta_t^A; \vartheta_t^B]$  satisfies the mean-flow equations

$$\frac{d}{dt}\vartheta_t = \begin{bmatrix} -R(\vartheta_t) & \gamma M^B(\vartheta_t) \\ \gamma M^A(\vartheta_t) & -R(\vartheta_t) \end{bmatrix} \vartheta_t - \begin{bmatrix} b(\vartheta_t) \\ b(\vartheta_t) \end{bmatrix}$$

with b defined in (41), where  $\theta = [\theta^A; \theta^B] \in \mathbb{R}^{2d}$  in the definition of  $\pi_{\theta}$ . The matrix R is unchanged from (46c), and

$$M^q(\theta) = \mathsf{E}_{\pi_{\theta}} [\psi_{(n)} \psi_{(n+1)}^{\theta^q}]^\intercal , \qquad q = A, B.$$

It is not clear how to establish stability of the mean flow using the techniques of this paper: the intuition following Prop. IV.2 is valid only if (following a transient) each of the policies  $\phi^{\vartheta_t^A}$ ,  $\phi^{\vartheta_t^B}$  approximates the  $\varepsilon$ -greedy policy for each of the two Q-function approximations.

**Zap Q-learning** Success requires that  $\partial_{\theta} \bar{f}(\theta)$  be non-singular for "most"  $\theta$ . Based on theory surrounding the Actor-Critic method, we have for a policy of the form (31),

$$\partial_{\theta} \bar{f}(\theta) = A_0(\theta) + Z(\theta), \qquad (51)$$

 $A_0(\theta) := \mathsf{E}_{\pi_{\theta}}[\partial_{\theta} f_{n+1}(\theta)] \text{ and } Z(\theta) := \mathsf{E}_{\pi_{\theta}}[\hat{f}_n(\theta)\Lambda_n(\theta)^{\intercal}].$  The random vectors in these definitions are:

- $f_{n+1}$  is given in (38).
- $\Lambda_n(\theta) = \nabla_{\theta} \log \widetilde{\Phi}^{\theta}(u \mid x)$ , evaluated at  $u = U_n$  and  $x = X_n$ , the score function associated with the randomized policy.
- $\hat{f}_n$  solves a certain Poisson equation. If the transition matrix  $T_{\theta}$  is aperiodic, then for a stationary realization of  $\{X_n, U_n : n \geq 0\}$  we have

$$\mathsf{E}_{\pi_{\theta}}[\hat{f}_{n}(\theta)\Lambda_{n}(\theta)^{\intercal}] = \sum_{k=0}^{\infty} \mathsf{E}_{\pi_{\theta}}[[f_{n-k}(\theta) - \bar{f}(\theta)]\Lambda_{n}(\theta)^{\intercal}]$$

Based on this representation we can obtain unbiased estimates of  $\partial_{\theta} \bar{f}(\theta_n)$  by adopting concepts from actor-critic algorithms. See [37, Ch. 10] for a survey in the style of this paper.

Analysis of the resulting algorithms will be considered in future research. The most likely path to success is to establish that the matrix norm of  $Z(\theta)$  is uniformly small, since approximations in this paper imply that  $A_0$  has desirable properties.

#### V. CONCLUSIONS

In the vast majority of application-oriented papers on reinforcement learning the policy used for training is not oblivious. Motivation for  $\varepsilon$ -greedy policies and their variants comes largely from the belief that optimism will lead to faster training. This paper makes clear that optimism is invaluable to ensure algorithmic stability.

There are of course a myriad of open questions.

- Can we obtain sharp bounds establishing that optimism leads to better sample complexity bounds (or perhaps better bounds on the asymptotic covariance)? The results in this paper combined with [15] might provide a first step. Ideally we want bounds sharp enough to inform the choice of  $\varepsilon$ , and for this it is likely best to begin with examination of asymptotic statistics.
- There is a large literature on SA with discontinuous dynamics, such as [6], [11]. It may be possible to extend Thm. IV.1 to  $\varepsilon$ -greedy policies (going beyond Thm. IV.5).
- Extension to average cost optimal control is straightforward through consideration of [1], as well as extension to relative Q-learning [21], [37].
- We should consider other paradigms for algorithm design. The recent approaches [34], [5], [30] are based on the linear programming formulation of optimal control, and are likely to lead to algorithms that respect desired performance bounds.

#### APPENDIX

This Appendix contains proofs of the main technical results concerning Q-learning subject to the linear function class assumption (5) and optimistic training. We begin with some general SA theory.

## A. Stochastic approximation theory

Convergence theory begins with a decomposition of the "disturbance" appearing in (10b):

$$\Delta_{n+1} = \mathcal{W}_{n+2} - \mathcal{T}_{n+2} + \mathcal{T}_{n+1} - \alpha_{n+1} \Upsilon_{n+2} \,. \tag{52}$$

Under the assumptions of [10], the dominating term in analysis is  $\{W_{n+2}\}$ , which is a martingale difference sequence. The sequence  $\{-\mathcal{T}_{n+2}+\mathcal{T}_{n+1}\}$  is telescoping so is insignificant in an ODE approximation. The sequence  $\{\alpha_{n+1}\Upsilon_{n+2}\}$  is small relative to the parameter sequence.

Lemma A.1 that follows provides bounds on the terms on the right hand side of (52) that are identical to those used in [10] to obtain the conclusions of Thm. II.1 and finer results. Consequently, the proof of Thm. II.1 is obtained by following identical steps in this prior work.

Moreover, if the limits in (24) hold, then the covariance matrix (25) admits the alternate representation

$$\Sigma_{\Lambda}^* = \mathsf{E}[\mathcal{W}_k^* \{ \mathcal{W}_k^* \}^{\mathsf{T}}]$$

in which  $\{\mathcal{W}_k^*: k \in \mathbb{Z}\}$  is a stationary version of the martingale difference sequence obtained from a stationary realization of  $\Phi^{\theta^*}$ .

The terms in (52) admit representations in terms of a family of solutions to Poisson's equation, following [36]. For each  $\theta$ , one version of the fundamental matrix associated with  $\mathcal{P}_{\theta}$  is the matrix inverse  $\mathcal{Z}_{\theta} = [I - \widetilde{\mathcal{P}}_{\theta}]^{-1}$ , where  $\widetilde{\mathcal{P}}_{\theta}(\xi, \xi') := \mathcal{P}_{\theta}(\xi, \xi') - \varpi_{\theta}(\xi')$  for each  $\xi, \xi'$ . Under (SA2) this may be expressed as the sum,  $\mathcal{Z}_{\theta} = \sum_{n=0}^{\infty} \widetilde{\mathcal{P}}_{\theta}^{n}$ . Writing  $\hat{f}(\theta, \xi) = \sum_{\xi'} \mathcal{Z}_{\theta}(\xi, \xi') f(\theta, \xi)$  and  $\Delta(\theta, \xi) := f(\theta, \xi) - \overline{f}(\theta)$ , the desired Poisson equation is solved:

$$\sum_{\xi'} \mathcal{P}_{\theta}(\xi, \xi') \hat{f}(\theta, \xi') = \hat{f}(\theta, \xi) - \Delta(\theta, \xi)$$

**Lemma A.1.** Under the assumptions of Thm. II.1,

$$\sup_{\theta} \left\| \partial_{\theta_{i}} \bar{f}\left(\theta\right) \right\| < \infty \quad \sup_{\xi, \theta} \left\| \partial_{\theta_{i}} \bar{f}\left(\theta, \xi\right) \right\| < \infty \tag{53}$$

Moreover, (52) holds with

$$\mathcal{W}_{n+2} := \hat{f}(\theta_n, \Phi_{n+2}) - \mathsf{E}[\hat{f}(\theta_n, \Phi_{n+2}) \mid \mathcal{F}_{n+1}], 
\mathcal{T}_{n+1} := \hat{f}(\theta_n, \Phi_{n+1}), 
\Upsilon_{n+2} := -\frac{1}{\alpha_{n+1}} [\hat{f}(\theta_{n+1}, \Phi_{n+2}) - \hat{f}(\theta_n, \Phi_{n+2})]$$
(54)

*Proof.* The decomposition with terms in (54) is obtained exactly as in [10, Section 2] following [36], so it remains to establish the bounds on  $\bar{f}$  and  $\hat{f}$  in (53).

Assumption (SA2) is a uniform Doeblin condition that implies  $\|\mathcal{Z}_{\theta}\|$  is uniformly bounded in  $\theta$  [39]. Write the invariance equation in operator theoretic notation  $\varpi_{\theta}\mathcal{P}_{\theta}=\varpi_{\theta}$ . On differentiating both sides we obtain the sensitivity formula of [44]:  $\partial_{\theta_i}\varpi_{\theta}=\varpi_{\theta}[\partial_{\theta_i}\mathcal{P}_{\theta}]\mathcal{Z}_{\theta}$ . Consequently, the invariant pmf  $\varpi_{\theta}$  enjoys the same smoothness properties as  $\mathcal{P}_{\theta}$ , giving  $\sup_{\xi,\theta,\underline{i}}|\partial_{\theta_i}\varpi_{\theta}(\xi)|\|\theta\|<\infty$ . This and (SA1) imply the bound on  $\partial_{\theta_i}f(\theta)$  in (53).

The bound for the solution to Poisson's equation follows from the identity  $\partial_{\theta_i} \mathcal{Z}_{\theta} = \mathcal{Z}_{\theta} \left[ \partial_{\theta_i} \widetilde{\mathcal{P}}_{\theta} \right] \mathcal{Z}_{\theta}$ . This combined with (SA1), (SA2) completes the proof of (53).

In the following we explain why the SA assumptions hold for Q-learning under the training policies of interest.

#### B. Validating SA assumptions for Q-learning

The proof of the following is postponed to the end of this subsection:

**Proposition A.2.** Under the assumptions of Thm. IV.1 the Q-learning algorithm satisfies Assumptions (SA1) and (SA2) with parameter-dependent noise  $\Phi_k = (X_k, I_k, D_k)$ .

Consider the oblivious policy defined by  $U_k \equiv W_k$  in the definition of  $R^{\mathcal{W}}$  in (46d). The transition matrix for the joint process  $\{(X_k, U_k) : k \geq 0\}$  can be obtained from (39b):

$$T_{\mathcal{W}}(z,z') = P_{\mathcal{U}}(x,x') \gamma_{\mathcal{W}}(u') ,$$

for any z=(x,u),  $z'=(x',u')\in X\times U$ . The invariance equation  $\pi_{\scriptscriptstyle{\mathcal{W}}}(z')=\sum_z\pi_{\scriptscriptstyle{\mathcal{W}}}(z)T_{\scriptscriptstyle{\mathcal{W}}}(z,z')$  implies that the invariant pmf is product form:

$$\pi_{\scriptscriptstyle \mathcal{W}}(z') = \mu_{\scriptscriptstyle \mathcal{W}}(x') \nu_{\scriptscriptstyle \mathcal{W}}(u')\,, \qquad z' = (x',u') \in \mathsf{X} \times \mathsf{U}\,.$$

in which  $\mu_{\mathcal{W}}(x') = \sum_{u} \pi_{\mathcal{W}}(x', u)$  is the steady-state marginal distribution of  $\boldsymbol{X}$  under this policy.

Similar notation is adopted for each of the invariant pmfs,

$$\mu_{\theta}(x) = \sum_{u} \pi_{\theta}(x, u), \qquad x \in X, \ \theta \in \mathbb{R}^{d}.$$

These are the invariant pmfs for the transition matrices  $\{P_{\theta}\}$ . Recall that these transition matrices and  $\{T_{\theta}\}$  are defined in (39).

Let  $X_0$  denote the support of  $\mu_{\mathcal{W}}$  and  $U_0$  the support of  $\nu_{\mathcal{W}}$ .

**Lemma A.3.** Suppose that (47a) holds, so that  $\pi_w$  is the unique invariant pmf. Consider any one of the three choices

of  $\{U_k\}$  used in (30) with  $\varepsilon < 1$  and any choice of  $\kappa$  in the case of (33) or  $\{\kappa_{\theta}\}$  in the case of (35). The following conclusions then hold:

(i)  $T_{\theta}$  is aperiodic, and for some  $N \geq 1$ ,  $\delta_T > 0$ ,

$$\sum_{k=1}^{N} T_{\theta}^{k}(z, z') \geq \delta_{T}, \quad z \in \mathsf{X} \times \mathsf{U}, \ z' \in \mathsf{X}_{0} \times \mathsf{U}_{0}, \ \theta \in \mathbb{R}^{d}.$$

- (ii) There is  $\delta_{\bullet} > 0$  such that  $\pi_{\theta}(z) \geq \delta_{\bullet} \pi_{w}(z)$  for all  $z, \theta$ .
- (iii)  $\pi_{\theta}(x, u) \geq \varepsilon \mu_{\theta}(x) \nu_{w}(u)$  for all  $x, u, \theta$ .
- (iv)  $R^{\mathcal{W}}(\theta) \geq \delta_{\bullet} R^{\mathcal{W}}$  for all  $\theta$ .

The constants  $\delta_T$ ,  $\delta_{\bullet}$  may depend on the policy parameters, but not  $\theta$ .

*Proof.* In view of (30) we have the bound  $T_{\theta}^{k}(z,z') \geq \varepsilon^{k}T_{\mathcal{W}}^{k}(z,z')$  for any k. Hence aperiodicity of  $T_{\theta}$  follows from the assumed aperiodicity of  $T_{\mathcal{W}}$ .

The lower bound in (i) holds for the oblivious policy under (47a): there is  $N \ge 1$  and  $\delta_w > 0$  such that

$$\sum_{k=1}^{N} T_{\mathcal{W}}^{k}(z, z') \ge \delta_{\mathcal{W}}, \quad for \ z \in \mathsf{X} \times \mathsf{U}, \ z' \in \mathsf{X}_{0} \times \mathsf{U}_{0}.$$

This implies the desired lower bound in (i) with  $\delta_T = \varepsilon^N \delta_w$ . Part (ii) follows from the bounds above and invariance:

$$\pi_{\theta}(z') = \sum_{z} \pi_{\theta}(z) \left(\frac{1}{N} \sum_{k=1}^{N} T_{\theta}^{k}(z, z')\right) \ge \frac{1}{N} \delta_{T}$$

Part (iii) also follows from invariance in the following onestep form: we have from (39b),

$$\begin{split} \pi_{\theta}(z') &= \sum_{z} \pi_{\theta}(z) T_{\theta}(z, z') \\ &= \sum_{x, u} \pi_{\theta}(x, u) P_{u}(x, x') \widetilde{\Phi}^{\theta}(u' \mid x') \\ &\geq \varepsilon \sum_{x, u} \pi_{\theta}(x, u) P_{u}(x, x') \nu_{w}(u') = \varepsilon \mu_{\theta}(x') \nu_{w}(u') \end{split}$$

The inequality follows from the bound  $\widetilde{\Phi}^{\theta}(u' \mid x') \geq \varepsilon \nu_{w}(u')$ . For part (iv), we begin with the definitions (46), giving

$$R^{\mathcal{W}}(\theta) = \sum_{x,u} \mu_{\theta}(x) \nu_{\mathcal{W}}(u) \psi(x,u) \psi(x,u)^{\mathsf{T}}$$

Applying (ii) gives  $\mu_{\theta}(x) \geq \delta_{\bullet} \mu_{w}(x)$  for all x, and hence the desired bound:

$$R^{\mathcal{W}}(\theta) \ge \delta_{\bullet} \sum_{x,u} \mu_{\mathcal{W}}(x) \nu_{\mathcal{W}}(u) \psi(x,u) \psi(x,u)^{\mathsf{T}} = \delta_{\bullet} R^{\mathcal{W}}$$

*Proof of Prop. A.2.* Assumption (SA1) is follows directly from the form of the recursion (6a), which gives  $f(\theta, \xi) = \psi(x, u)[c(x, u) + \gamma \min_{u'} \{\psi(x', u')^{\mathsf{T}}\theta\}]$ , in which x, u, x' is a function of  $\xi = (x; \iota; \delta)$ .

Part (i) of (SA2) follows directly from Lemma A.3. It remains to establish (ii).

For this we apply (15) and (39a), which imply it is enough to show that a similar bound holds for the policy:

$$\sup_{u,x,\theta} \left| \partial_{\theta_i} \widetilde{\Phi}^{\theta}(u \mid x) \right| \|\theta\| < \infty$$

This follows from the construction, giving  $\widetilde{\Phi}^{\theta} = \widetilde{\Phi}^{r\theta}$  for each  $r \geq 1$  and parameter satisfying  $\|\theta\| \geq 1$ .

## C. Mean flow for the $\varepsilon$ -greedy policy

In this subsection the input is chosen to be the  $\varepsilon$ -greedy policy (30). The motivation is in part the fact that establishing stability of the ODE@ $\infty$  in this case is far easier than the tamed Gibbs approximation.

The transition matrix (39b) becomes

$$T_{\theta}(z, z') = P_{u}(x, x') \left\{ (1 - \varepsilon) \mathbb{1} \left\{ u' = \phi^{\theta}(x') \right\} + \varepsilon \nu^{\mathcal{W}}(u') \right\}$$
 (55)

for  $z=(x,u), \ z'=(x',u')\in X\times U$ . The family  $\{T_{\theta}:\theta\in\mathbb{R}^d\}$  is finite because there are only a finite number of deterministic stationary policies; it takes on a constant value on each connected component of  $\mathcal{C}^{\Theta}$  (recall (29)).

Compact representations of f and f are obtained with additional notation. For  $n \ge 0$  denote

$$\psi_{(n)}^{\Theta} = \psi(X_n, \phi^{\theta_n}(X_n)) \qquad \psi_{(n)}^{W} = \psi(X_n, \mathcal{W}_n)$$

$$c_n^{\Theta} = c(X_n, \phi^{\theta_n}(X_n)) \qquad c_n^{W} = c(X_n, \mathcal{W}_n)$$
(56)

We have under the  $\varepsilon$ -greedy policy (30, 32),

$$f_{n+1}(\theta_n) = (1 - B_n) \left( c_n^{\Theta} + \left[ \gamma \psi_{(n+1)}^{\Theta} - \psi_{(n)}^{\Theta} \right]^{\mathsf{T}} \theta_n \right) \psi_{(n)}^{\Theta}$$
$$+ B_n \left( c_n^{W} + \left[ \gamma \psi_{(n+1)}^{\Theta} - \psi_n^{W} \right]^{\mathsf{T}} \theta_n \right) \psi_{(n)}^{W}$$
 (57)

**Lemma A.4.**  $\mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} \}^{\intercal} \right] = R_{-1}(\theta) + \varepsilon D(\theta)$ , in which

$$D(\theta) = \mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} - \psi_{(n+1)}^{\mathcal{W}} \}^{\mathsf{T}} \right] \tag{58}$$

Proof. Starting with the definition

$$R_{-1}(\theta) = \mathsf{E}_{\pi_{\theta}} \big[ \psi_{(n)} \{ \psi_{(n+1)} \}^{\mathsf{T}} \big]$$

we have under the  $\varepsilon$ -greedy policy,  $R_{-1}(\theta) =$ 

$$\begin{split} &(1-\varepsilon)\mathsf{E}_{\pi_{\theta}}\big[\psi_{(n)}\{\psi_{(n+1)}^{\Theta}\}^{\intercal}\big] + \varepsilon\mathsf{E}_{\pi_{\theta}}\big[\psi_{(n)}\{\psi_{(n+1)}^{\mathcal{W}}\}^{\intercal}\big] \\ &= \mathsf{E}_{\pi_{\theta}}\big[\psi_{(n)}\{\psi_{(n+1)}^{\Theta}\}^{\intercal}\big] + \varepsilon\mathsf{E}_{\pi_{\theta}}\big[\psi_{(n)}\{\psi_{(n+1)}^{\mathcal{W}} - \psi_{(n+1)}^{\Theta}\}^{\intercal}\big] \end{split}$$

**Lemma A.5.** The vector fields for the mean flow and the  $ODE@\infty$  for the  $\varepsilon$ -greedy policy are

$$\bar{f}(\theta) = A(\theta)\theta - b(\theta)$$
  $\bar{f}_{\infty}(\theta) = A(\theta)\theta$  (59a)

in which 
$$A(\theta) = -[R_0(\theta) - \gamma R_{-1}(\theta)] + \varepsilon \gamma D(\theta)$$
 (59b)  
 $b(\theta) = (1 - \varepsilon)b^{\Theta}(\theta) + \varepsilon b^{W}(\theta)$  (59c)

$$b^{\scriptscriptstyle \Theta}(\theta) = -\mathsf{E}_{\pi_{\theta}}[\psi_{(n)}^{\scriptscriptstyle \Theta}c_n^{\scriptscriptstyle \Theta}] \ \ \text{and} \ \ b^{\scriptscriptstyle \mathcal{W}}(\theta) = -\mathsf{E}_{\pi_{\theta}}[\psi_{(n)}^{\scriptscriptstyle \mathcal{W}}c(X_n,\mathcal{W}_n)].$$

*Proof.* The representation (57) is equivalently expressed  $f_{n+1}(\theta_n) = A_{n+1}\theta_n - b_{n+1}$ , in which

$$A_{n+1} = \psi_{(n)} \left[ \gamma \psi_{(n+1)}^{\Theta} - \psi_{(n)} \right]^{\mathsf{T}} b_{n+1} = (1 - B_n) \psi_{(n)}^{\Theta} c(X_n, \Phi^{\theta}(X_n)) + B_n \psi_{(n)}^{\mathcal{W}} c(X_n, \mathcal{W}_n)$$

The expression for  $b(\theta)$  in the expression  $\bar{f}(\theta) = \mathsf{E}_{\pi_{\theta}}[f_{n+1}(\theta)] = A(\theta)\theta - b(\theta)$  is immediate.

We have  $A(\theta) = -R_0(\theta) + \gamma \mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} \}^{\mathsf{T}} \right]$ , so that (59b) follows from Lemma A.4.

The expression for  $\bar{f}_{\infty}$  follows from the fact that A and b are invariant under positive scaling of their arguments:  $A(r\theta) = A(\theta)$  and  $b(r\theta) = b(\theta)$  for any  $\theta$  and r > 0.

The mean flow (11) is a differential inclusion because the vector field  $\bar{f}$  is not continuous.

The form of the expression for  $A(\theta)$  in (59b) is intended to evoke the similar formula (IV.3) obtained for (49).

The following conclusions are based on arguments similar to what is used to obtain stability of on-policy TD-learning [50]. Recall the definition (48):  $\varepsilon_{\gamma} := (1-\gamma)^2/[(1-\gamma)^2+\gamma^2]$ .

**Proposition A.6.** If  $\varepsilon < \varepsilon_{\gamma}$ , then there is  $\beta_{\varepsilon} > 0$  such that  $v^{\mathsf{T}}A(\theta)v \leq -\beta_{\varepsilon}||v||^2$  for each  $v, \theta \in \mathbb{R}^d$ .

*Proof.* Applying Lemma A.5 gives for any  $v, \theta$ ,

$$v^{\mathsf{T}} A(\theta) v \le -(1 - \gamma) v^{\mathsf{T}} R_0(\theta) v + \varepsilon \gamma v^{\mathsf{T}} D(\theta) v \tag{60}$$

The inequality follows from the bound  $v^{\mathsf{T}}R_k(\theta)v \leq v^{\mathsf{T}}R_0(\theta)v$ , valid for any k.

We are left to bound the term involving D. Write  $v^{\mathsf{T}}D(\theta)v = d_v^{\Theta}(\theta) - d_v^{W}(\theta)$  with

$$\begin{split} d_v^{\Theta}(\theta) &= \mathsf{E}_{\pi_{\theta}} \big[ (v^{\mathsf{T}} \psi_{(n)}) (v^{\mathsf{T}} \psi_{(n+1)}^{\Theta}) \big] \\ d_v^{\mathcal{W}}(\theta) &= \mathsf{E}_{\pi_{\theta}} \big[ (v^{\mathsf{T}} \psi_{(n)}) (v^{\mathsf{T}} \psi_{(n+1)}^{\mathcal{W}}) \big] \end{split}$$

Using the bound  $xy \leq \frac{1}{2}[x^2 + y^2]$  for  $x, y \in \mathbb{R}$ , we obtain for any  $\delta_{\mathcal{W}}, \delta_{\Theta} > 0$ ,

$$\begin{aligned} \left| d_v^{\Theta}(\theta) \right| &\leq \frac{1}{2} \delta_{\Theta}^{-1} v^{\mathsf{T}} R_0(\theta) v + \frac{1}{2} \delta_{\Theta} v^{\mathsf{T}} R_0^{\Theta}(\theta) v \\ \left| d_v^{W}(\theta) \right| &\leq \frac{1}{2} \delta_{\mathcal{W}}^{-1} v^{\mathsf{T}} R_0(\theta) v + \frac{1}{2} \delta_{\mathcal{W}} v^{\mathsf{T}} R_0^{W}(\theta) v \end{aligned}$$

Recall from (46c) that  $R_0(\theta) = (1 - \varepsilon)R_0^{\Theta}(\theta) + \varepsilon R_0^{W}(\theta)$ . Set  $\delta_{W} = \varepsilon \eta$ ,  $\delta_{\Theta} = (1 - \varepsilon)\eta$ , with  $\eta > 0$  to be chosen. Then,

$$\begin{split} v^{\mathsf{T}}D(\theta)v &\leq \frac{1}{2} \Big[ \left( \delta_{\mathcal{W}}^{-1} + \delta_{\Theta}^{-1} \right) v^{\mathsf{T}} R_0(\theta) v \\ &+ \delta_{\Theta} v^{\mathsf{T}} R_0^{\Theta}(\theta) v + \delta_{\mathcal{W}} v^{\mathsf{T}} R_0^{\mathcal{W}}(\theta) v \Big] \\ &= \frac{1}{2} \Big[ \left( \frac{1}{\varepsilon} + \frac{1}{1 - \varepsilon} \right) \frac{1}{\eta} + \eta \Big] v^{\mathsf{T}} R_0(\theta) v \end{split}$$

The value  $\eta_{\varepsilon}^* = \sqrt{\varepsilon^{-1} + (1 - \varepsilon)^{-1}}$  minimizes the right hand side, and on substitution,  $v^{\mathsf{T}}D(\theta)v \leq \eta_{\varepsilon}^* v^{\mathsf{T}}R_0(\theta)v$ .

Substitution into (60) gives the final bound,

$$\boldsymbol{v}^{\mathsf{T}}\boldsymbol{A}(\boldsymbol{\theta})\boldsymbol{v} \leq \left[-(1-\gamma) + \varepsilon\gamma\eta_{\varepsilon}^{*}\right]\boldsymbol{v}^{\mathsf{T}}\boldsymbol{R}_{0}(\boldsymbol{\theta})\boldsymbol{v}$$

The coefficient is negative for positive  $\varepsilon$  if and only if  $\varepsilon < \varepsilon_{\gamma}$ . We obtain the desired bound with

$$\beta_{\varepsilon} = \left[ (1 - \gamma) - \varepsilon \gamma \eta_{\varepsilon}^{*} \right] \min_{\theta} \lambda_{\min}(R_{0}(\theta))$$

Lemma A.3 implies that the minimum is strictly positive.

The extension of Prop. A.6 to the tamed Gibbs policy requires approximations summarized in the next subsection.

#### D. Entropy and Gibbs bounds

Consider a single Gibbs pmf on U with energy  $E \colon \mathsf{U} \to \mathbb{R}$  and inverse temperature  $\kappa > 0$ :

$$p_{\kappa}(u) = \frac{1}{\mathcal{Z}_{\kappa}} \exp(-\kappa E(u)), \quad u \in \mathsf{U},$$

The normalizing factor  $\mathcal{Z}_{\kappa}$  is known as the *partition function*. The entropy of  $p_{\kappa}$  is denoted

$$H_{\kappa} = -\sum_{u} p_{\kappa}(u) \log(p_{\kappa}(u)) = \sum_{u} p_{\kappa}(u) \left[ \kappa E(u) + \log(\mathcal{Z}_{\kappa}) \right]$$

Bounds on entropy imply bounds on the quality of the softmin approximation. Denote  $E := \min_{u} E(u)$ .

**Lemma A.7.** For any  $\kappa > 0$ ,

$$\underline{E} \le \sum_{u} p_{\kappa}(u) E(u) \le \underline{E} + \frac{1}{\kappa} \log(|\mathsf{U}|)$$

Proof. The uniform distribution maximizes entropy, giving

$$\sum_{u} p_{\kappa}(u) \left[ \kappa E(u) + \log(\mathcal{Z}_{\kappa}) \right] \le \log(|\mathsf{U}|)$$

The bound  $\log(\mathcal{Z}_\kappa)=\log\sum_u\exp(-\kappa E(u))\geq -\kappa\underline{E}$  completes the proof.

An implication of the lemma to the policy (35): for any initial distribution for  $(X_0, U_0)$ ,

$$\underline{Q}^{\theta}(X_{k+1}) \leq \mathsf{E}\left[Q^{\theta}(X_{k+1}, \mathcal{U}_{k+1}) \mid X_0^{k+1}, U_0^k\right] 
\leq \underline{Q}^{\theta}(X_{k+1}) + \frac{1}{\kappa_{\theta}} \log(|\mathsf{U}|), \quad k \geq 0.$$
(61)

## E. Proof of Thms. IV.1 and IV.5

The proof of Thm. IV.1 closely follows the proof of Prop. A.6. We begin a companion to Lemma A.4:

**Lemma A.8.** We have for the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy,

$$\mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} \}^{\mathsf{T}} \right] = R_{-1}(\theta) + \varepsilon D(\theta) + (1 - \varepsilon) E(\theta) \tag{62a}$$

in which 
$$D(\theta) = \mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} - \psi_{(n+1)}^{\mathcal{W}} \}^{\mathsf{T}} \right]$$
 (62b)  
 $E(\theta) = \mathsf{E}_{\pi_{\theta}} \left[ \psi_{(n)} \{ \psi_{(n+1)}^{\Theta} - \psi_{(n+1)}^{\mathcal{U}} \}^{\mathsf{T}} \right]$  (62c)

with 
$$\psi_{(n+1)}^{\mathcal{U}} = \psi(X_{n+1}, \mathcal{U}_{n+1}).$$

We have a partial extension of Prop. A.6:

**Lemma A.9.** The following holds for the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy, subject to (47) and  $\varepsilon < \varepsilon_{\gamma}$ : there is  $\beta_{\varepsilon} > 0$  such that  $\theta^{\intercal}A(\theta)\theta \leq -\beta_{\varepsilon}\|\theta\|^2$  for all  $\kappa_0 > 0$  sufficiently large, and all  $\|\theta\| \geq 1$ .

*Proof.* Applying Lemma A.8 to (41), and following the same steps as in the proof of Prop. A.6 we obtain

$$\begin{split} \theta^{\mathsf{T}} A(\theta) \theta &\leq -\beta_{\varepsilon}^{0} \, \theta^{\mathsf{T}} R_{0}(\theta) \theta + \gamma (1-\varepsilon) \theta^{\mathsf{T}} E(\theta) \theta \\ \textit{with} \quad \beta_{\varepsilon}^{0} &= \left[ (1-\gamma) - \varepsilon \gamma \eta_{\varepsilon}^{*} \right] > 0 \\ \eta_{\varepsilon}^{*} &= \sqrt{\varepsilon^{-1} + (1-\varepsilon)^{-1}} \end{split}$$

From the definition (62c) we have

$$\theta^\intercal E(\theta)\theta = \mathsf{E}_{\pi_\theta} \big[ Q^\theta(X_n, U_n) \{ \underline{Q}^\theta(X_{n+1}) - Q^\theta(X_{n+1}, \mathcal{U}_{n+1}) \} \big]$$

Applying (61) and the expression for  $\kappa_{\theta}$  in (36), we obtain for  $\|\theta\| \ge 1$ ,

$$\begin{aligned} \left| \theta^{\mathsf{T}} E(\theta) \theta \right| &\leq \frac{1}{\kappa_0} \|\theta\| \log(|\mathsf{U}|) \mathsf{E}_{\pi_{\theta}} \left[ \left| Q^{\theta}(X_n, U_n) \right| \right] \\ &\leq \frac{1}{\kappa_0} \|\theta\|^2 \log(|\mathsf{U}|) \sqrt{\lambda_{\text{max}}} \end{aligned}$$

with  $\lambda_{\max}$  the maximum over all  $\theta$  of the maximum eigenvalue of  $R_0(\theta)$ . Combining these bounds completes the proof.

*Proof of Thm. IV.1.* Precisely as in the proof of Prop. IV.2 we obtain a solution to (v4) using  $V(\theta) = \|\theta\|$  (recall (21)), which implies (23) exactly as in the case when  $\Phi$  is exogenous.

The existence of  $\theta^*$  follows from Prop. II.2, exactly as in the proof in Prop. IV.2

*Proof of Thm. IV.5.* Let  $\theta^{\kappa_0}$  denote the solution to the projected Bellman equation for the  $(\varepsilon, \kappa_0)$ -tamed Gibbs policy, in which  $\varepsilon < \varepsilon_\gamma$  is fixed.

Observe that in Lemma A.9 we obtain a uniform bound over all large  $\kappa_0$ . An examination of the proof of Prop. II.2 shows that there is a constant  $b_{\varepsilon}$  such that  $\|\theta^{\kappa_0}\| \leq b_{\varepsilon}$  for all sufficiently large  $\kappa_0$ .

Hence we can find a subsequence  $\kappa_0^n \to \infty$  as  $n \to \infty$ , for which the following limits exist:

$$\theta^* = \lim_{n \to \infty} \theta^{\kappa_0^n}, \quad \pi^* = \lim_{n \to \infty} \pi_n,$$

in which  $\pi_n$  is the invariant pmf obtained from the policy using  $\theta^{\kappa_0^n}$ .

The invariant pmfs have the form

$$\pi_n(x, u) = \mu_n(x)\widetilde{\phi}^n(u \mid x)$$

with  $\widetilde{\Phi}^n$  defined in (35) using  $\kappa_0^n$ , and  $\mu_n$  the first marginal of  $\pi_n$ . It follows that the limiting invariant pmf has the same structure,  $\pi^*(x,u) = \mu^*(x)\widetilde{\Phi}^{\theta^*}(u\mid x)$ . Since  $\kappa_0^n\uparrow\infty$ , convergence implies that  $\widetilde{\Phi}^{\theta^*}$  is of the form (50) with  $\widetilde{\Phi}^*\in\Phi^{\theta^*}$ .

Letting  $\bar{f}_n$  denote the vector field obtained using  $\theta^{\kappa_0^n}$  we must have convergence for each  $\theta$ :

$$\bar{f}(\theta) = \lim_{n \to \infty} \bar{f}_n(\theta) = \mathsf{E}_{\pi^*}[\psi_{(n)}\mathcal{B}(X_n, U_n; \theta)],$$

in which  $U_n$  is defined using the randomized  $\varepsilon$ -greedy policy  $\widetilde{\Phi}^{\theta^*}$ , and  $\mathcal{B}$  defined in (40b) is a continuous function of  $\theta$ . Since  $\overline{f}_n(\theta_n)=0$  for each n, we conclude that  $\overline{f}(\theta^*)=0$  as desired.

## REFERENCES

- J. Abounadi, D. Bertsekas, and V. S. Borkar. Learning algorithms for Markov decision processes with average cost. SIAM Journal on Control and Optimization, 40(3):681–698, 2001.
- [2] S. Asmussen and P. W. Glynn. Stochastic Simulation: Algorithms and Analysis, volume 57 of Stochastic Modelling and Applied Probability. Springer-Verlag, New York, 2007.
- [3] K. E. Avrachenkov, V. S. Borkar, H. P. Dolhare, and K. Patil. Full gradient DQN reinforcement learning: A provably convergent scheme. In *Modern Trends in Controlled Stochastic Processes*:, pages 192–220. Springer, 2021.
- [4] L. Baird. Residual algorithms: Reinforcement learning with function approximation. In A. Prieditis and S. Russell, editors, *Proc. Machine Learning*, pages 30–37. Morgan Kaufmann, San Francisco (CA), 1995.
- [5] J. Bas Serrano, S. Curi, A. Krause, and G. Neu. Logistic Q-learning. In A. Banerjee and K. Fukumizu, editors, *Proc. of The Intl. Conference on Artificial Intelligence and Statistics*, volume 130, pages 3610–3618, 13–15 Apr 2021.
- [6] M. Benaïm, J. Hofbauer, and S. Sorin. Stochastic approximations and differential inclusions, Part II: applications. *Mathematics of Operations Research*, 31(4):673–695, 2006.
- [7] A. Benveniste, M. Métivier, and P. Priouret. Adaptive algorithms and stochastic approximations, volume 22. Springer Science & Business Media, Berlin Heidelberg, 2012.
- [8] D. P. Bertsekas. Dynamic programming and optimal control. Vol. II. Athena Scientific, Belmont, MA, fourth edition, 2012.
- [9] S. Bhatnagar. The Borkar–Meyn Theorem for asynchronous stochastic approximations. Systems & control letters, 60(7):472–478, 2011.

- [10] V. Borkar, S. Chen, A. Devraj, I. Kontoyiannis, and S. Meyn. The ODE method for asymptotic statistics in stochastic approximation and reinforcement learning. arXiv e-prints:2110.14427, pages 1–50, 2021.
- [11] V. S. Borkar. Stochastic approximation with 'controlled Markov' noise. Systems & control letters, 55(2):139–145, 2006.
- [12] V. S. Borkar. Stochastic Approximation: A Dynamical Systems Viewpoint. Hindustan Book Agency, Delhi, India, 2nd edition, 2021.
- [13] V. S. Borkar and S. P. Meyn. The ODE method for convergence of stochastic approximation and reinforcement learning. SIAM J. Control Optim., 38(2):447–469, 2000.
- [14] S. Chen, A. M. Devraj, F. Lu, A. Bušić, and S. Meyn. Zap Q-Learning with nonlinear function approximation. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Proc. Conference on Neural Information Processing Systems (NeurIPS), and arXiv e-prints* 1910.05405, volume 33, pages 16879–16890, 2020.
- [15] Z. Chen, J.-P. Clarke, and S. T. Maguluri. Target network and truncation overcome the deadly triad in Q-learning. SIAM Journal on Mathematics of Data Science, 5(4):1078–1101, 2023.
- [16] A. Cooper and S. Meyn. Reinforcement learning design for quickest change detection. Submitted for publication, and arXiv preprint arXiv:2403.14109, 2024.
- [17] D. De Farias and B. Van Roy. On the existence of fixed points for approximate value iteration and temporal-difference learning. *Journal* of Optimization Theory and Applications, 105(3):589–608, 2000.
- [18] A. M. Devraj. Reinforcement Learning Design with Optimal Learning Rate. PhD thesis, University of Florida, 2019.
- [19] A. M. Devraj, A. Bušić, and S. Meyn. Fundamental design principles for reinforcement learning algorithms. In K. G. Vamvoudakis, Y. Wan, F. L. Lewis, and D. Cansever, editors, *Handbook on Reinforcement Learning* and Control, Studies in Systems, Decision and Control series (SSDC, volume 325). Springer, 2021.
- [20] A. M. Devraj and S. P. Meyn. Zap Q-learning. In Proc. of the Intl. Conference on Neural Information Processing Systems, pages 2232– 2241, 2017.
- [21] A. M. Devraj and S. P. Meyn. Q-learning with uniformly bounded variance. *IEEE Trans. on Automatic Control*, 67(11):5948–5963, 2022.
- [22] A. Gopalan and G. Thoppe. Approximate Q-learning and SARSA(0) under the ε-greedy policy: a differential inclusion analysis. arXiv preprint arXiv:2205.13617, 2022.
- [23] H. V. Hasselt. Double Q-learning. In Proc. Advances in Neural Information Processing Systems, pages 2613–2621, 2010.
- [24] T. Jaakola, M. Jordan, and S. Singh. On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation*, 6:1185–1201, 1994.
- [25] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143, 2020.
- [26] P. Karmakar and S. Bhatnagar. Stochastic approximation with iterate-dependent Markov noise under verifiable conditions in compact state space with the stability of iterates not ensured. *IEEE Transactions on Automatic Control*, 66(12):5941–5954, 2021.
- [27] C. K. Lauand and S. Meyn. Revisiting step-size assumptions in stochastic approximation. arXiv 2405.17834, 2024.
- [28] D. Lee and N. He. A unified switching system perspective and ODE analysis of Q-learning algorithms. arXiv, page arXiv:1912.02270, 2019.
- [29] H.-D. Lim, D. W. Kim, and D. Lee. Regularized Q-learning. preprint arXiv:2202.05404, 2022.
- [30] F. Lu, P. G. Mehta, S. P. Meyn, and G. Neu. Convex analytic theory for convex Q-learning. In *IEEE Conference on Decision and Control*, pages 4065–4071, Dec 2022.
- [31] H. R. Maei, C. Szepesvári, S. Bhatnagar, and R. S. Sutton. Toward off-policy learning control with function approximation. In *Proc. ICML*, pages 719–726, USA, 2010. Omnipress.
- [32] A. Martinelli, M. Gargiani, M. Draskovic, and J. Lygeros. Data-driven optimal control of affine systems: A linear programming perspective. *IEEE Control Systems Letters*, 6:3092–3097, 2022.
- [33] A. Martinelli, M. Gargiani, and J. Lygeros. Data-driven optimal control with a relaxed linear program. *Automatica*, 136:110052, 2022.
- [34] P. G. Mehta and S. P. Meyn. Q-learning and Pontryagin's minimum principle. In *Proc. of the Conf. on Dec. and Control*, pages 3598–3605, Dec. 2009.
- [35] F. S. Melo, S. P. Meyn, and M. I. Ribeiro. An analysis of reinforcement learning with function approximation. In *Proc. ICML*, pages 664–671, New York, NY, 2008.
- [36] M. Metivier and P. Priouret. Theoremes de convergence presque sure pour une classe d'algorithmes stochastiques a pas decroissants. *Prob. Theory Related Fields*, 74:403–428, 1987.

- [37] S. Meyn. Control Systems and Reinforcement Learning. Cambridge University Press, Cambridge, 2022.
- [38] S. Meyn. Stability of Q-learning through design and optimism. *arXiv* 2307.02632, 2023.
- [39] S. P. Meyn and R. L. Tweedie. Markov chains and stochastic stability. Cambridge University Press, Cambridge, second edition, 2009. Published in the Cambridge Mathematical Library. 1993 edition online.
- [40] B. T. Polyak. A new method of stochastic approximation type. Avtomatika i telemekhanika (in Russian). translated in Automat. Remote Control, 51 (1991), pages 98–107, 1990.
- [41] A. Ramaswamy and S. Bhatnagar. A generalization of the Borkar-Meyn Theorem for stochastic recursive inclusions. *Mathematics of Operations Research*, 42(3):648–661, 2017.
- [42] G. A. Rummery and M. Niranjan. On-line Q-learning using connectionist systems. Technical report 166, Cambridge Univ., Dept. Eng., Cambridge, U.K. CUED/F-INENG/, 1994.
- [43] D. Ruppert. Efficient estimators from a slowly convergent Robbins-Monro processes. Technical Report Tech. Rept. No. 781, Cornell University, School of Operations Research and Industrial Engineering, Ithaca, NY, 1988.
- [44] P. J. Schweitzer. Perturbation theory and finite Markov chains. J. Appl. Prob., 5:401–403, 1968.
- [45] R. Sutton and A. Barto. Reinforcement Learning: An Introduction. MIT Press, Cambridge, MA, 2nd edition, 2018.
- [46] R. S. Sutton. Learning to predict by the methods of temporal differences. *Mach. Learn.*, 3(1):9–44, 1988.
- [47] C. Szepesvári. Algorithms for Reinforcement Learning. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2010.
- [48] C. Szepesvari, E. Brunskill, S. Bubeck, A. Malek, S. Meyn, A. Tewari, and M. Wang. Theory of Reinforcement Learning Boot Camp. Aug 31 to Sep 4, 2020. https://simons.berkeley.edu/workshops/rl-2020-bc.
- [49] J. Tsitsiklis. Asynchronous stochastic approximation and Q-learning. Machine Learning, 16:185–202, 1994.
- [50] J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Trans. Automat. Control*, 42(5):674–690, 1997.
- [51] B. Van Roy. Learning and Value Function Approximation in Complex Decision Processes. PhD thesis, Massachusetts Institute of Technology, Cambridge, MA, 1998. AAI0599623.
- [52] M. J. Wainwright. Stochastic approximation with cone-contractive operators: Sharp  $\ell_{\infty}$ -bounds for *Q*-learning. *CoRR*, abs/1905.06265, 2019.
- [53] C. J. C. H. Watkins. Learning from Delayed Rewards. PhD thesis, King's College, Cambridge, Cambridge, UK, 1989.
- [54] C. J. C. H. Watkins and P. Dayan. Q-learning. Machine Learning, 8(3-4):279–292, 1992.
- [55] L. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756, 2020.
- [56] F. Zarin Faizal and V. Borkar. Functional Central Limit Theorem for two timescale stochastic approximation. arXiv e-prints, page arXiv:2306.05723, June 2023.



Sean Meyn was raised by the beach in Santa Monica, California. Following his BA in mathematics at UCLA, he moved on to pursue a PhD with Peter Caines at McGill University. After about 20 years as a professor of ECE at the University of Illinois, in 2012 he moved to beautiful Gainesville. He is now Professor and Robert C. Pittman Eminent Scholar Chair in the Department of Electrical and Computer Engineering at the University of Florida, director of the Laboratory for Cognition and Control, and Inria International Chair at Inria, France. He is an IEEE

CSS distinguished lecturer. His interests span many aspects of stochastic control, stochastic processes, information theory, and optimization. For the past decade, his applied research has focused on engineering, markets, and policy in energy systems.