

A Appendices

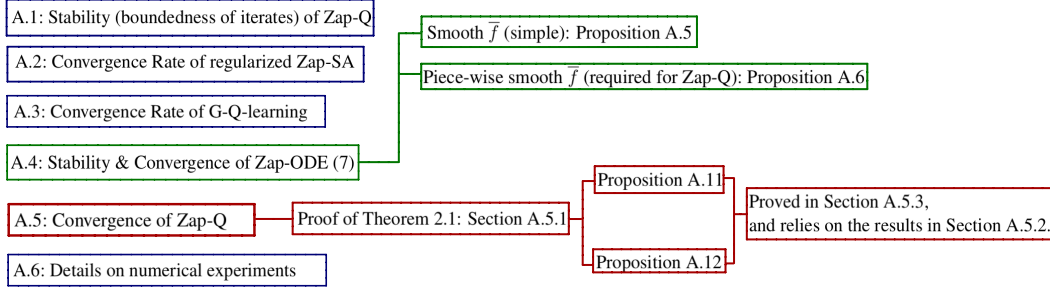


Figure 2: Organization of appendices.

Figure 2 contains an overview of the appendices to help the reader navigate through the technical results leading to the proof of the main results, summarized in Thm. 2.1.

Sections A.1–A.3 collect proofs of miscellaneous technical claims: Section A.1 contains sufficient conditions to ensure boundedness of $\{\theta_n\}$. Section A.2 is focused on the linearization of Zap SA algorithm (14), and based on this a formula for the asymptotic covariance of Zap SA is obtained. The proof that the maximal eigenvalue for GQ-learning with tabular basis satisfies $\text{Real}(\lambda(A_*) \geq -(1 - \gamma)^2$ is contained in Section A.3.

Analysis of the Zap ODE (7) is contained in Section A.4, including existence of solutions, and the consistency result $\lim_{t \rightarrow \infty} \bar{f}(w_t) = 0$. Section A.5 finishes the proof of Thm. 2.1 by establishing solidarity between the ODE and the stochastic recursion. It is simplest to begin in the special case in which ζ has non-negative entries, which is the focus of Sections A.5.2 and A.5.3. The arguments are extended in Section A.5.4, where the non-negativity is relaxed.

A.1 Establishing boundedness of the parameter estimates

Suppose that the following limits exist:

$$Q_\infty^\theta(x, u) = \lim_{m \rightarrow \infty} m^{-1} Q^{m\theta}(x, u), \quad x \in \mathcal{X}, u \in \mathcal{U}$$

$$\zeta_\infty(\theta, z) = \lim_{m \rightarrow \infty} \zeta(m\theta, z), \quad z \in \mathcal{Z}$$

where the limiting functions are twice continuously differentiable. The global Lipschitz conditions in (A2) imply that the gradients also converge, and the convergence is uniform on compact sets. We then obtain a vector field for the “ODE at infinity” introduced in [10]:

$$\bar{f}_\infty(\theta) := \lim_{m \rightarrow \infty} m^{-1} \bar{f}(m\theta) = \mathbb{E}[(\gamma Q_\infty^\theta(X_{n+1}) - Q_\infty^\theta(X_n, U_n)) \zeta_\infty(\theta, \Phi_n)]$$

and a similar definition for $f_\infty(\theta, z)$. The associated regularized Newton-Raphson flow “at infinity” is similar to (7):

$$\frac{d}{dt} w_t = -[\varepsilon I + A_\infty(w_t)^\top A_\infty(w_t)]^{-1} A_\infty(w_t)^\top \bar{f}_\infty(w_t), \quad A_\infty(w_t) = \partial_\theta \bar{f}_\infty(w_t) \quad (32)$$

With $\mathcal{A}_\infty(\theta)$ defined as in (26) with respect to \bar{f}_∞ , assume the following:

(A2_∞) The functions Q_∞ and ζ_∞ are Lipschitz continuous and twice continuously differentiable in θ in any open set not containing the origin; $f_\infty(\theta, z)$ is Lipschitz continuous for each $z \in \mathcal{Z}$; $A^\top \bar{f}_\infty(\theta) \neq 0$ for all $\theta \neq 0$ and $A \in \mathcal{A}_\infty(\theta)$.

Assumptions (A2) and (A2_∞) are identical when the function approximation Q^θ is linear, and $\zeta = \nabla Q^\theta$.

The function $\|\bar{f}_\infty\|$ is coercive under (A2_∞) since \bar{f}_∞ is radially linear: $\bar{f}_\infty(m\theta) = m\bar{f}_\infty(\theta)$ for any θ and any $m \geq 0$. Prop. A.6 can be adapted to show that (32) is globally asymptotically stable. [9, Sections 6.3, Theorem 9] explains how stability of the ODE implies stability of the SA algorithm.

A.2 Asymptotic covariance of regularized Zap SA

We first introduce a standard result in linear system theory [24, Theorem 2.6-1].

Lemma A.1. *If $A \in \mathbb{R}^{d \times d}$ is Hurwitz and $\Sigma_\Delta \in \mathbb{R}^{d \times d}$ is positive semi-definite, then there exists a unique solution $\Sigma \geq 0$ that solves the Lyapunov equation,*

$$A\Sigma + \Sigma A^\top + \Sigma_\Delta = 0,$$

whose solution can be expressed

$$\Sigma = \int_0^\infty \exp(A\tau) \Sigma_\Delta \exp(A^\top \tau) d\tau$$

Let $\{\mathcal{E}_n^G\}$ be the sequence obtained by the stochastic linear recursion (10) with matrix gain $G \in \mathbb{R}^{d \times d}$:

$$\mathcal{E}_{n+1}^G = \mathcal{E}_n^G + \alpha_{n+1} G[A_* \mathcal{E}_n^G + \Delta_{n+1}], \quad \mathcal{E}_0^G = \theta_0 - \theta^* \quad (33)$$

Denote the asymptotic covariance of $\{\mathcal{E}_n^G\}$ by $\Sigma_\theta^G := \lim_{n \rightarrow \infty} n \mathbb{E}[\mathcal{E}_n^G (\mathcal{E}_n^G)^\top]$. According to the eigenvalue test (12), Σ_θ^G is finite if $\frac{1}{2}I + GA_*$ is Hurwitz. It is well known that the matrix gain $G_* = -A_*^{-1}$ achieves the minimal asymptotic covariance:

$$\Sigma_\theta^* = A_*^{-1} \Sigma_\Delta (A_*^{-1})^\top \quad (34)$$

The following result is standard in stochastic approximation [5, Part I, Section 3.2.3, Proposition 4], and quantifies optimality of Σ_θ^* .

Lemma A.2. *Suppose that A_* and $\frac{1}{2}I + GA_*$ are Hurwitz. Then,*

(i) *The asymptotic covariance $\Sigma_\theta^G \geq 0$ uniquely solves the Lyapunov equation:*

$$(\frac{1}{2}I + GA_*) \Sigma_\theta^G + \Sigma_\theta^G (\frac{1}{2}I + GA_*)^\top + G \Sigma_\Delta G^\top = 0$$

(ii) *The sub-optimality gap $\tilde{\Sigma}_\theta^G = \Sigma_\theta^G - \Sigma_\theta^* \geq 0$ uniquely solves the Lyapunov equation:*

$$(\frac{1}{2}I + GA_*) \tilde{\Sigma}_\theta^G + \tilde{\Sigma}_\theta^G (\frac{1}{2}I + GA_*)^\top + (G + A_*^{-1}) \Sigma_\Delta (G + A_*^{-1})^\top = 0 \quad (35)$$

For any symmetric matrix $S \in \mathbb{R}^{d \times d}$, denote by $\lambda(S)$ the set of its eigenvalues.

Proposition A.3. *Suppose $A_* \in \mathbb{R}^{d \times d}$ is Hurwitz, and denote $G_\varepsilon = -[\varepsilon I + A_*^\top A_*]^{-1} A_*^\top$. If $0 < \varepsilon < \lambda_{\min}(A_*^\top A_*)$, then $\frac{1}{2}I + G_\varepsilon A_*$ is Hurwitz, so that the matrix gain G_ε in the linear recursion (33) results in a finite asymptotic covariance $\Sigma_\theta^\varepsilon$. Moreover, the follow approximation holds:*

$$\Sigma_\theta^\varepsilon = \Sigma_\theta^* + \varepsilon^2 \Sigma_\theta^{(2)} + O(\varepsilon^3), \quad \text{with } \Sigma_\theta^{(2)} = (A_* A_*^\top A_*)^{-1} \Sigma_\Delta (A_*^\top A_* A_*^\top)^{-1}. \quad (36)$$

Proof. The set of eigenvalues of $\frac{1}{2}I + G_\varepsilon A_*$ admits the following representations:

$$\begin{aligned} \lambda(\frac{1}{2}I + G_\varepsilon A_*) &= \left\{ \frac{1}{2} - \lambda : \lambda \in \lambda([\varepsilon I + A_*^\top A_*]^{-1} A_*^\top A_*) \right\} \\ &= \left\{ \frac{1}{2} - \frac{1}{\lambda} : \lambda \in \lambda(\varepsilon(A_*^\top A_*)^{-1} + I) \right\} \\ &= \left\{ \frac{1}{2} - \frac{1}{\varepsilon\lambda + 1} : \lambda \in \lambda((A_*^\top A_*)^{-1}) \right\} \\ &= \left\{ \frac{1}{2} - \frac{1}{\varepsilon/\lambda + 1} : \lambda \in \lambda(A_*^\top A_*) \right\} \end{aligned}$$

Given $0 < \varepsilon < \lambda_{\min}(A_*^\top A_*)$, the eigenvalues of $\frac{1}{2}I + G_\varepsilon A_*$ are real and strictly negative. In particular, this matrix is Hurwitz, as claimed.

We next establish the approximation (36). By Lemma A.2 (iii), $\tilde{\Sigma}_\theta^\varepsilon = \Sigma_\theta^\varepsilon - \Sigma_\theta^*$ solves the Lyapunov equation (35) with G replaced by G_ε . Denoting $\tilde{G}_\varepsilon = G_\varepsilon + A_*^{-1}$, we obtain by Lemma A.1,

$$\tilde{\Sigma}_\theta^\varepsilon = \int_0^\infty \exp([\frac{1}{2}I + G_\varepsilon A_*]\tau) \tilde{G}_\varepsilon \Sigma_\Delta \tilde{G}_\varepsilon^\top \exp([\frac{1}{2}I + G_\varepsilon A_*]^\top \tau) d\tau \quad (37)$$

A Taylor series representation of matrix inverse results in the following:

$$\begin{aligned} G_\varepsilon A_* &= -[\varepsilon I + A_*^\top A_*]^{-1} A_*^\top A_* = -[I + \varepsilon(A_*^\top A_*)^{-1}]^{-1} = -[I - \varepsilon A_*^{-1}(A_*^\top)^{-1}] + O(\varepsilon^2) \\ \tilde{G}_\varepsilon &= G_\varepsilon + A_*^{-1} = [G_\varepsilon A_* + I] A_*^{-1} = \varepsilon A_*^{-1}(A_*^\top)^{-1} A_*^{-1} + O(\varepsilon^2) \end{aligned}$$

With $\Sigma_\theta^{(2)} = (A_* A_*^\top A_*)^{-1} \Sigma_\Delta (A_*^\top A_* A_*^\top)^{-1}$, the integral in (37) becomes

$$\tilde{\Sigma}_\theta^\varepsilon = \varepsilon^2 \int_0^\infty \exp(-[\tfrac{1}{2}I - \varepsilon A_*^{-1}(A_*^\top)^{-1}]\tau) \Sigma_\theta^{(2)} \exp(-[\tfrac{1}{2}I - \varepsilon A_*^{-1}(A_*^\top)^{-1}]^\top \tau) d\tau + O(\varepsilon^3) \quad (38)$$

Another Taylor series expansion for matrix exponential gives

$$\exp(-[\tfrac{1}{2}I - \varepsilon A_*^{-1}(A_*^\top)^{-1}]\tau) = \exp(-\tfrac{\tau}{2}I) + \varepsilon \tau A_*^{-1}(A_*^\top)^{-1} \exp(-\tfrac{\tau}{2}I) + O(\varepsilon^2)$$

Consequently, the integral in (38) can be rewritten as

$$\begin{aligned} \tilde{\Sigma}_\theta^\varepsilon &= \varepsilon^2 \int_0^\infty \exp(-\tfrac{\tau}{2}I) \Sigma_\theta^{(2)} \exp(-\tfrac{\tau}{2}I) d\tau + O(\varepsilon^3) \\ &= \varepsilon^2 \Sigma_\theta^{(2)} + O(\varepsilon^3) \end{aligned}$$

□

A.3 Eigenvalue test for GQ-learning

Consider the linear function approximation architecture: $Q^\theta(x, u) = \psi(x, u)^\top \theta$, where $\psi : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}^d$ is the basis function. With eligibility vector $\zeta_n := \psi(X_n, U_n)$, let \bar{f} be defined by (18). GQ-learning [31] aims to solve the root finding problem (18), transformed into the non-convex optimization problem (3).

The GQ-learning [31] algorithm is the two-time scale SA algorithm,

$$\theta_{n+1} = \theta_n + \alpha_{n+1} [\mathcal{D}(\theta_n, \Phi_{n+1}) \zeta_n - \gamma \varphi_{n+1}^\top \zeta_n \psi(X_{n+1}, \phi^{\theta_n}(X_{n+1}))] \quad (39a)$$

$$\varphi_{n+1} = \varphi_n + \beta_{n+1} \zeta_n [\mathcal{D}(\theta_n, \Phi_{n+1}) - \psi(X_n, U_n)^\top \varphi_n] \quad (39b)$$

where $\{\beta_n\}$ and $\{\alpha_n\}$ are non-negative step-size sequences satisfying $\alpha_n/\beta_n \rightarrow 0$ as $n \rightarrow \infty$. The fast time scale recursion (39b) for $\{\varphi_n\}$ is designed so that $\varphi_n \approx M \bar{f}(\theta_n)$ for large n . It follows that the ODE approximation of (39a) is (4).

Proposition A.4. *The linearization matrix for GQ-learning at θ^* is given by $A_{GQ} = -A_*^\top M A_*$, whenever θ^* is a solution to $A(\theta)^\top M \bar{f}(\theta) = 0$. With the tabular basis: $\psi_k(x, u) = \mathbb{I}\{x = x^k, u = u^k\}$, $1 \leq k \leq \ell_x \cdot \ell_u$, there is an eigenvalue λ_{GQ} of A_{GQ} satisfying*

$$\lambda_{GQ} \geq -(1 - \gamma)^2$$

We first introduce some notation for tabular Q-learning. For any deterministic stationary policy $\phi : \mathcal{X} \rightarrow \mathcal{U}$, let S_ϕ denote the substitution operator, defined for any function $Q : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ by $S_\phi Q(x) = Q(x, \phi(x))$. With P viewed as a matrix with $\ell_x \cdot \ell_u$ rows and ℓ_x columns, PS_ϕ can be interpreted as the transition matrix for the joint process (\mathbf{X}, \mathbf{U}) when \mathbf{U} is defined using policy ϕ [15]. Then $\bar{f}(\theta)$ can be written in matrix form

$$\bar{f}(\theta) = \Pi r + \Pi[\gamma PS_{\phi^\theta} - I]\theta \quad (40)$$

where Π is a diagonal matrix with entries: $\Pi(k, k) := \varpi(x^k, u^k)$ and r is a vector with entries: $r(k) := r(x^k, u^k)$. The derivative $A(\theta)$ of $\bar{f}(\theta)$ is given by

$$A(\theta) = \Pi[\gamma PS_{\phi^\theta} - I]$$

Proof. The matrix A_{GQ} is the derivative of $-A(w_t)^\top M \bar{f}(w_t)$ at θ^* . For the tabular case, by (40),

$$A_{GQ} = -[\gamma PS_{\phi^*} - I]^\top \Pi [\gamma PS_{\phi^*} - I] = -H^\top H,$$

with $H := \Pi^{1/2}[I - \gamma PS_{\phi^*}]$. It suffices to show that $H^\top H$ has a positive eigenvalue less than $(1 - \gamma)^2$.

Since H^{-1} is a positive and irreducible matrix, we can apply the same arguments as in [15, Theorem 3.3] to bound the Perron-Frobenius eigenvalue as follows:

$$\lambda_{\text{PF}} \geq \frac{1}{1-\gamma} \min_{x,u} \frac{1}{\sqrt{\varpi(x,u)}}$$

Therefore, H has positive eigenvalue $\lambda_H = \lambda_{\text{PF}}^{-1}$ such that

$$\lambda_H \leq (1-\gamma) \max_{x,u} \sqrt{\varpi(x,u)}$$

Applying [21, Theorem 5.6.9] we obtain the complementary bound

$$\lambda_H \geq \sigma_{\min}(H) = \sqrt{\lambda_{\min}(H^\top H)}$$

and combining the two implies:

$$\lambda_{\min}(H^\top H) \leq \lambda_H^2 \leq (1-\gamma)^2 \left(\max_{x,u} \sqrt{\varpi(x,u)} \right)^2 \leq (1-\gamma)^2$$

□

A.4 ODE Analysis

To obtain the existence of a solution to (7), we first consider an ideal smooth setting:

Proposition A.5. *Consider the following conditions for the function \bar{f} :*

- (a) *\bar{f} is globally Lipschitz continuous and continuously differentiable. Hence $A(\cdot)$ is a bounded matrix-valued function.*
- (b) *$\|\bar{f}\|$ is coercive. That is, $\{\theta : \|\bar{f}(\theta)\| \leq n\}$ is compact for each n .*
- (c) *The function \bar{f} has a unique zero θ^* , and $A^\top(\theta)\bar{f}(\theta) \neq 0$ for $\theta \neq \theta^*$. Moreover, the matrix $A_* = A(\theta^*)$ is non-singular.*

The following hold for solutions to the ODE (7) under increasingly stronger assumptions:

- (i) *If (a) holds then for each t , and each initial condition*

$$\frac{d}{dt} \bar{f}(w_t) = -A(w_t)[\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t) \quad (41)$$

- (ii) *If in addition (b) holds, then the solutions to the ODE are bounded, and*

$$\lim_{t \rightarrow \infty} A(w_t)^\top \bar{f}(w_t) = 0 \quad (42)$$

- (iii) *If (a)–(c) hold, then (7) is globally asymptotically stable.*

□

Proof. The result (i) follows from the chain rule and the definitions.

The proof of (ii) is based on the Lyapunov function $V(w) = \frac{1}{2} \|\bar{f}(w)\|^2$ combined with (a):

$$\frac{d}{dt} V(w_t) = -\bar{f}(w_t)^\top A(w_t)[\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t)$$

The right hand side is non-positive when $w_t \neq \theta^*$. Integrating each side gives for any $T > 0$,

$$V(w_T) = V(w_0) - \int_0^T \bar{f}(w_t)^\top A(w_t)[\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t) dt \quad (43)$$

so that $V(w_T) \leq V(w_0)$ for all T . Under the coercive assumption, it follows that solutions to (7) are bounded. Also, letting $T \rightarrow \infty$, we obtain from (43) the bound

$$\int_0^\infty \bar{f}(w_t)^\top A(w_t)[\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t) dt \leq V(w_0)$$

This combined with boundedness of w_t implies that $\lim_{t \rightarrow \infty} A(w_t)^\top \bar{f}(w_t) = 0$.

We next prove (iii). Global asymptotic stability of (7) requires that solutions converge to θ^* from each initial condition, and also that θ^* is stable in the sense of Lyapunov [26]. Assumption (c) combined with (ii) gives the former, that $\lim_{t \rightarrow \infty} w_t = \theta^*$. A convenient sufficient condition for the latter is obtained by considering $A_\infty = \partial_\theta [\mathcal{G}(\theta) \bar{f}(\theta)]|_{\theta=\theta^*}$. Stability in the sense of Lyapunov holds if this matrix is Hurwitz (all eigenvalues are in the strict left half plane in \mathbb{C}) [26, Thm. 4.7]. Apply the definitions, we obtain $A_\infty = -[\varepsilon I + M]^{-1} M$ with $M = A(\theta^*)^\top A(\theta^*) > 0$ (recall that $A(\theta^*)$ is assumed to be non-singular). The matrix A_∞ is negative definite, and hence Hurwitz.

□

Prop. A.5 cannot be applied to the ODE (7) that motivated Zap Q-learning because \bar{f} is only piecewise smooth. To obtain an extension we consider the ODE in its integral form:

$$w_t = w_0 - \int_0^t [\varepsilon I + A^\top(w_\tau) A(w_\tau)]^{-1} A^\top(w_\tau) \bar{f}(w_\tau) d\tau, \quad t \geq 0 \quad (44)$$

where $\bar{f}(\theta)$, $A(\theta)$ are defined in (18, 20).

Proposition A.6. *Under Assumptions A1-A2, there exists a solution to (44) from each initial condition. The following hold for any solution:*

- (i) $\bar{f}(w_t) = \bar{f}(w_0) - \int_0^t A(w_\tau) [\varepsilon I + A(w_\tau)^\top A(w_\tau)]^{-1} A(w_\tau)^\top \bar{f}(w_\tau) d\tau, \quad t \geq 0$
- (ii) $\|\bar{f}(w_t)\|$ is non-increasing, and $\lim_{t \rightarrow \infty} \bar{f}(w_t) = 0$.
- (iii) If in addition A3 holds, then the ODE (44) is globally asymptotically stable.

The proof of existence is obtained by considering smooth approximations of (7).

Lemma A.7. *Under Assumptions A1-A2, there exists a solution to (44) from each initial condition.*

Proof. Define a C^∞ probability density η on \mathbb{R}^d via

$$\eta(x) := \begin{cases} k \exp(-(1 - \|x\|^2)^{-1}) & \|x\| < 1, \\ 0 & \|x\| \geq 1, \end{cases} \quad (45)$$

where $k > 0$ is a normalization constant: $\int \eta(x) dx = 1$. For each $\delta > 0$, a C^∞ vector field is defined via the convolution:

$$\bar{f}_\delta(x) = \frac{1}{\delta^d} \int \bar{f}(x - y) \eta(y/\delta) dy, \quad x \in \mathbb{R}^d \quad (46)$$

The family of functions $\{\bar{f}_\delta : 0 < \delta \leq 1\}$ is globally uniformly Lipschitz continuous, with the same Lipschitz constant b_L as of \bar{f} . It is also evident that $\lim_{\delta \downarrow 0} \bar{f}_\delta = \bar{f}$ pointwise. The uniform Lipschitz continuity implies that the convergence is uniform on compact sets.

Denote $A_\delta(\theta) = \partial_\theta \bar{f}_\delta(\theta)$, and consider the ODE (44) with \bar{f} and A replaced by their smooth approximations:

$$w_t^\delta = w_0^\delta - \int_0^t [\varepsilon I + A_\delta^\top(w_\tau^\delta) A_\delta(w_\tau^\delta)]^{-1} A_\delta^\top(w_\tau^\delta) \bar{f}_\delta(w_\tau^\delta) d\tau, \quad w_0^\delta = w_0 \quad (47)$$

The solution exists and is unique for each $\delta \in (0, 1]$. To obtain bounds on the solution we require bounds on the matrices involved, and opt for the spectral norm:

$$\begin{aligned} \|\varepsilon I + A_\delta^\top(w_t^\delta) A_\delta(w_t^\delta)\|^{-1} &= \frac{1}{\lambda_{\min}(\varepsilon I + A_\delta^\top(w_t^\delta) A_\delta(w_t^\delta))} \leq \frac{1}{\varepsilon} \\ \|A_\delta(w_\tau^\delta)\| &\leq b_L \end{aligned}$$

where b_L is the Lipschitz constant for \bar{f}_δ . Therefore,

$$\begin{aligned}
\|w_t^\delta\| &\leq \|w_0^\delta\| + \int_0^t \|[\varepsilon I + A_\delta^\top(w_\tau^\delta)A_\delta(w_\tau^\delta)]^{-1}\| \cdot \|A_\delta(w_\tau^\delta)\| \cdot \|\bar{f}_\delta(w_\tau^\delta)\| d\tau \\
&\leq \|w_0^\delta\| + \frac{b_L}{\varepsilon} \int_0^t \|\bar{f}_\delta(w_\tau^\delta)\| d\tau \\
&\leq \|w_0^\delta\| + \frac{b_L}{\varepsilon} \int_0^t \|\bar{f}_\delta(w_0^\delta)\| + \|\bar{f}_\delta(w_\tau^\delta) - \bar{f}_\delta(w_0^\delta)\| d\tau \\
&\leq \|w_0^\delta\| + \frac{b_L}{\varepsilon} \{T\|\bar{f}_\delta(w_0^\delta)\| + b_L \int_0^t \|w_\tau^\delta - w_0^\delta\| d\tau\} \\
&\leq \|w_0^\delta\| + \frac{b_L}{\varepsilon} \{T(\|\bar{f}_\delta(w_0^\delta)\| + b_L\|w_0^\delta\|) + b_L \int_0^t \|w_\tau^\delta\| d\tau\}
\end{aligned}$$

The set $\{\|\bar{f}_\delta(w_0^\delta)\| : 0 < \delta \leq 1\}$ is bounded by $\max_{y \in \mathcal{B}(w_0, 1)} \|\bar{f}(y)\|$, where $\mathcal{B}(w_0, 1)$ denotes the closed unit ball in \mathbb{R}^d centered at w_0 . By Gronwall's inequality, there exist constants C_1 and C_2 such that

$$\|w_t^\delta\| \leq C_1 + C_2 e^{b_L^2 T / \varepsilon}, \quad t \in [0, T], \quad \delta \in (0, 1]$$

This combined with (47) implies that $\{w^\delta : 0 < \delta \leq 1\}$ is uniformly bounded and equicontinuous. By the Arzelà-Ascoli theorem, there exists a sequence $\delta_n \downarrow 0$ and a continuous function $w^0 : [0, T] \rightarrow \mathbb{R}^d$ such that

$$\lim_{n \rightarrow \infty} \sup_{t \in [0, T]} \|w_t^{\delta_n} - w_t^0\| = 0$$

So the functional equation (47) holds for w^0 with $\delta = 0$, and w^0 is thus a solution of (44). \square

The following result has been derived in [15, Lemma A.10]. We present it here for completeness.

Lemma A.8. *Let $G(\theta) := \max_{1 \leq i \leq \ell_u} G_i(\theta)$ where each $G_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is twice continuously differentiable and Lipschitz continuous. Let $w : [0, T] \rightarrow \mathbb{R}^d$ be a Lipschitz continuous function, and denote $g_t := G(w_t)$. Then,*

- (i) $g : [0, T] \rightarrow \mathbb{R}$ is Lipschitz continuous.
- (ii) At any time $t_0 \in (0, T)$ such that the derivatives of g_t and w_t exist,

$$\left. \frac{d}{dt} g_t \right|_{t=t_0} = \partial_\theta G_k(w_{t_0}) \cdot \left. \frac{d}{dt} w_t \right|_{t=t_0} \quad \text{for each } k \in \arg \max_i G_i(w_{t_0}). \quad (48)$$

Proof. Denote $g_t^i = G_i(w_t)$, so that $g_t = \max_{1 \leq i \leq \ell_u} g_t^i$. Let b_L denote a Lipschitz constant for each of these functions:

$$|g_{t_1}^i - g_{t_0}^i| \leq b_L |t_1 - t_0|, \quad t_0, t_1 \in [0, T], \quad 1 \leq i \leq \ell_u$$

For any $t_0, t_1 \in [0, T]$,

$$\begin{aligned}
g_{t_1} - g_{t_0} &\leq g_{t_1}^k - g_{t_0}^k, \quad \text{for each } k \in \arg \max_i g_{t_0}^i \\
&\leq b_L |t_1 - t_0|
\end{aligned}$$

The same inequality holds for $g_{t_0} - g_{t_1}$ with $k \in \arg \max_i g_{t_1}^i$. This proves (i).

The proof of (ii) is also straightforward: The difference $g_t - g_t^k$ has a global minimum at t_0 if $k \in \arg \max_i g_{t_0}^i$, and consequently

$$0 = \left. \frac{d}{dt} [g_t - g_t^k] \right|_{t=t_0}$$

\square

Given a parameter vector $\theta \in \mathbb{R}^d$, denote by $\varsigma^\theta : \mathbf{X} \times \mathbf{U} \rightarrow \mathbb{R}$ the reward function that satisfies the Bellman equation (16), with Q^* replaced by Q^θ : For each $x \in \mathbf{X}$ and $u \in \mathbf{U}$,

$$\varsigma^\theta(x, u) := -\gamma \sum_{x' \in \mathbf{X}} P_u(x, x') \underline{Q}^\theta(x') + Q^\theta(x, u) \quad (49)$$

Lemma A.9. *Suppose Assumptions A1-A2 hold and the function $w : [0, T] \rightarrow \mathbb{R}^d$ is Lipschitz continuous. Then, $\varsigma^{w_t}(x, u)$ is Lipschitz continuous in t for each x, u . Moreover, at any point t_0 of differentiability,*

$$\left. \frac{d}{dt} \varsigma^{w_t}(x, u) \right|_{t=t_0} = \left[-\gamma \sum_{x' \in \mathbf{X}} P_u(x, x') \partial_\theta Q^{w_{t_0}}(x', \phi^{w_{t_0}}(x')) + \partial_\theta Q^{w_{t_0}}(x, u) \right] \left. \frac{d}{dt} w_t \right|_{t=t_0} \quad (50)$$

where $\phi^{w_{t_0}}$ is defined in (17).

Proof. From the definition (49), it is sufficient to establish the derivative formula

$$\left. \frac{d}{dt} \underline{Q}^{w_t}(x') \right|_{t=t_0} = \partial_\theta Q^{w_{t_0}}(x', \phi^{(k)}(x')) \cdot \left. \frac{d}{dt} w_t \right|_{t=t_0}$$

where $\phi^{(k)}$ is any policy that is $Q^{w_{t_0}}$ -greedy. This is immediate from Lemma A.8. \square

Stability of is obtained from the following standard Lyapunov condition:

Lemma A.10. *Suppose that $\{w_t : t \in \mathbb{R}\}$ is a Lipschitz continuous function taking values in \mathbb{R}^d , $V : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is continuous and coercive, and $U : \mathbb{R}^d \rightarrow \mathbb{R}_+$. Assume moreover the following properties:*

- (i) $\inf\{U(\theta) : V(\theta) \geq \delta\} > 0$ for each $\delta > 0$.
- (ii) $V(w_t) \leq V(w_0) - \int_0^t U(w_\tau) d\tau, \quad t \geq 0$.

Then, there exists a function $B : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $V(w_t) \leq \eta$ for all $t \geq V(w_0)B(\eta)$. In particular, $V(w_t) \rightarrow 0$ as $t \rightarrow \infty$.

Proof. For any scalar η satisfying $0 < \eta < V(w_0)$, let $H^\eta := \{\theta : \eta \leq \theta \leq V(w_0)\}$ and

$$\varepsilon_\eta = \inf_{\theta \in H^\eta} U(\theta)$$

Under assumption (i) of the lemma we have $\varepsilon_\eta > 0$. Let $T^\eta = \inf\{t : V(w_t) \leq \eta\}$, so that $w_t \in H^\eta$ for $0 \leq t \leq T^\eta$. By assumption (ii) we have

$$0 \leq V(w_t) \leq V(w_0) - \varepsilon_\eta t, \quad 0 < t \leq T^\eta.$$

Therefore, $T^\eta < V(w_0)/\varepsilon_\eta$. Because $V(w_t)$ is non-increasing in t , we have $V(w_t) \leq \eta$ for all $t \geq V(w_0)B(\varepsilon)$, with $B(\varepsilon) = \varepsilon_\eta^{-1}$.

Since η is arbitrary, it follows that $\lim_{t \rightarrow \infty} V(w_t) = 0$. \square

Proof of Prop. A.6. Suppose $w : [0, T] \rightarrow \mathbb{R}^d$ is a solution of (44). At point t of differentiability, the derivative of $\bar{f}(w_t)$ is given by

$$\begin{aligned} \frac{d}{dt} \bar{f}(w_t) &= \frac{d}{dt} \mathbb{E} \left[\zeta_n \varsigma^{w_t}(X_n, U_n) \right] \\ &= \mathbb{E} \left[\zeta_n \frac{d}{dt} \varsigma^{w_t}(X_n, U_n) \right] + \mathbb{E} \left[\mathcal{D}(w_t, \Phi_{n+1}) \frac{d}{dt} \zeta_n \right] \end{aligned} \quad (51)$$

For each $x \in \mathbf{X}$ and $u \in \mathbf{U}$, $\varsigma^{w_t}(x, u)$ is a Lipschitz continuous function of t , whose derivative is given in Lemma A.9. Assertion (i) follows:

$$\begin{aligned} \frac{d}{dt} \bar{f}(w_t) &= \mathbb{E} \left[\zeta_n [\gamma \partial_\theta Q^{w_t}(X_{n+1}, \phi^{w_t}(X_{n+1})) - \partial_\theta Q^{w_t}(X_n, U_n)] + \mathcal{D}(w_t, \Phi_{n+1}) \partial_\theta \zeta_n \right] \frac{d}{dt} w_t \\ &= -A(w_t) [\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t) \end{aligned}$$

A candidate Lyapunov function is defined as $V(w_t) := \frac{1}{2} \|\bar{f}(w_t)\|^2$. At a point t where $\bar{f}(w_t)$ is differentiable,

$$\begin{aligned} \frac{d}{dt} V(w_t) &= -\bar{f}(w_t)^\top A(w_t) [\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top \bar{f}(w_t) \\ &\quad - [\varepsilon I + A(\theta)^\top A(\theta)]^{-1} \leq -b_V I \end{aligned} \quad (52)$$

The integral representation of (52) then gives, for any $t \in [0, T]$,

$$\begin{aligned} V(w_t) &= V(w_0) - \int_0^t \bar{f}(w_\tau)^\top A(w_\tau) [\varepsilon I + A(w_\tau)^\top A(w_\tau)]^{-1} A(w_\tau)^\top \bar{f}(w_\tau) d\tau \\ &\leq V(w_0) - b_V \int_0^t \|A(w_\tau)^\top \bar{f}(w_\tau)\|^2 d\tau \end{aligned} \quad (53)$$

Under (A2) the assumptions of Lemma A.10 hold with $U(\theta) = b_V \|A(\theta)^\top \bar{f}(\theta)\|^2$, so that $\lim_{t \rightarrow \infty} V(w_t) = \lim_{t \rightarrow \infty} \bar{f}(w_t) = 0$.

If in addition (A3) holds, we conclude that $\lim_{t \rightarrow \infty} w_t = \theta^*$. \square

A.5 Proof of Thm. 2.1

The remainder of the Appendix is dedicated to the proof of Thm. 2.1. We use $n_0 = 0$ in the definition of the step-size sequences (22); this shortens many of the expressions that follow, and the extension to general $n_0 \geq 1$ is obvious. Given the typical choice of ζ_n in (2), it is assumed throughout that $\zeta_n := \zeta(\theta_n, X_n, U_n)$ for some function $\zeta : \mathbb{R}^d \times \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}^d$. We proceed under the additional assumption that the vector-valued function ζ has non-negative entries:

$$[\zeta(\theta, x, u)]_i \geq 0 \text{ for each } i, \theta, x, u. \quad (54)$$

The proofs are extended to the general case in Section A.5.4.

A.5.1 Generalities

This subsection contains the building blocks of the proof, summarized in two propositions, and the proof of Thm. 2.1 based on these key results. The proofs of the propositions are postponed to subsequent subsections.

The *slow time scale* used for an ODE approximation of $\{\theta_n\}$ is defined by

$$t_n = \sum_{i=1}^n \alpha_i = \sum_{i=1}^n \frac{1}{i}, \quad n \geq 1, \quad t_0 = 0 \quad (55)$$

and its approximate inverse

$$[t] := \max\{j : t_j \leq t\} \quad (56)$$

Define the continuous time process $\{\bar{w}_t : t \geq 0\}$ with $\bar{w}_{t_n} = \theta_n$, and extended to \mathbb{R}_+ via linear interpolation. Define the associated continuous time process $\{\bar{c}_t := \bar{f}(\bar{w}_t) : t \geq 0\}$. We also define the piecewise constant processes $\{\bar{\mathcal{A}}_t, \bar{\mathcal{G}}_t : t \geq 0\}$ with $\bar{\mathcal{A}}_t = \hat{A}_{n+1}, \bar{\mathcal{G}}_t = G_{n+1}$ for $t \in [t_n, t_{n+1})$. Both $b_\theta := \sup_n \|\theta_n\| = \sup_t \|\bar{w}_t\|$ and $b_c := \sup_t \|\bar{c}_t\|$ are finite *a.s.* by assumption.

Denote by $\phi(1) = \phi(T_0, t)$ a function of two variables, satisfying for each $T > 0$,

$$\lim_{T_0 \rightarrow \infty} \sup_{0 \leq t \leq T} \|\phi(T_0, t)\| = 0$$

Proposition A.11. *Under Assumptions (A1)-(A2) and (54), $\{\bar{w}_t\}$ and $\{\bar{c}_t\}$ are Lipschitz continuous with respect to t , and the following approximations hold:*

- (i) $\lim_{T_0 \rightarrow \infty} \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t \frac{d}{dt} \bar{w}_t - \frac{d}{dt} \bar{c}_t\|_\infty dt = 0.$
- (ii) $\bar{c}_{T_0+t} = \bar{c}_{T_0} + \int_{T_0}^{T_0+t} \bar{\mathcal{A}}_\tau \bar{\mathcal{G}}_\tau \bar{c}_\tau d\tau + \phi(1), \quad T_0 \rightarrow \infty$

$$(iii) \quad \|\bar{c}_{T_0+t}\|^2 = \|\bar{c}_{T_0}\|^2 + 2 \int_{T_0}^{T_0+t} \bar{c}_\tau^\top \bar{\mathcal{A}}_\tau \bar{\mathcal{G}}_\tau \bar{c}_\tau d\tau + o(1), \quad T_0 \rightarrow \infty. \quad \square$$

For a fixed but arbitrary time-horizon $T > 0$, define a family of functions $\{\bar{\Gamma}^{T_0} : T_0 \geq 0\}$, where $\bar{\Gamma}^{T_0} : [0, T] \rightarrow \mathbb{R}^m$ for each $T_0 \geq 0$ and an integer m . It consists of four components: for $t \in [0, T]$,

$$\bar{\Gamma}_1^{T_0}(t) = \bar{w}_{T_0+t}, \quad \bar{\Gamma}_2^{T_0}(t) = \bar{c}_{T_0+t}, \quad \bar{\Gamma}_3^{T_0}(t) = \bar{\mathcal{A}}_{T_0+t}, \quad \bar{\Gamma}_4^{T_0}(t) = -\bar{\mathcal{A}}_{T_0+t} \bar{\mathcal{G}}_{T_0+t}$$

$\{\bar{\Gamma}_1^{T_0} : T_0 \geq 0\}$ and $\{\bar{\Gamma}_2^{T_0} : T_0 \geq 0\}$ are uniformly Lipschitz continuous and bounded. More specifically, each of $\bar{\Gamma}_1^{T_0}$ and $\bar{\Gamma}_2^{T_0}$ is a function of two variables: $\bar{\Gamma}_1^{T_0}(\omega, t), \bar{\Gamma}_2^{T_0}(\omega, t)$ with $\omega \in \Omega$ and $t \in \mathbb{R}_+$. The property that $\bar{\Gamma}_1^{T_0}$ and $\bar{\Gamma}_2^{T_0}$ are Lipschitz continuous and bounded holds with probability one. Denote their sub-sequential limits by

$$\Gamma_1(t) = w_t, \quad \Gamma_2(t) = c_t$$

where the convergence is uniform over $[0, T]$.

Limits of the remaining components of Γ are defined with respect to the weak topology in $L_2([0, T]; \mathbb{R}^{d \times d})$. Because $\{\bar{\Gamma}_3^{T_0} : T_0 \geq 0\}$ and $\{\bar{\Gamma}_4^{T_0} : T_0 \geq 0\}$ are uniformly bounded, they are weakly relatively sequentially compact in $L_2([0, T]; \mathbb{R}^{d \times d})$ [17, Theorem 1.1.2]. Their weak sub-sequential limits Γ_3 and Γ_4 are denoted by $\{\mathcal{A}_t, \mathcal{H}_t : 0 \leq t \leq T\}$. That is, there exists $T_k \rightarrow \infty$ such that

$$\bar{\Gamma}_3^{T_k} \rightarrow \mathcal{A} \text{ weakly in } L_2([0, T]; \mathbb{R}^{d \times d}), \quad \bar{\Gamma}_4^{T_k} \rightarrow \mathcal{H} \text{ weakly in } L_2([0, T]; \mathbb{R}^{d \times d}), \quad k \rightarrow \infty$$

Based on Prop. A.11, and a separate analysis of the *fast time scale* recursion for $\{\hat{A}_n\}$ we obtain the following properties for any sub-sequential limit Γ of $\{\bar{\Gamma}^{T_0} : T_0 \geq 0\}$:

Proposition A.12. *Under Assumptions (A1)-(A2) and (54), for each $t \in [0, T]$,*

- (i) $c_t := \Gamma_2(t) = \bar{f}(w_t)$.
- (ii) $\mathcal{A}_t := \Gamma_3(t) \in \mathcal{A}(w_t)$.
- (iii) $\mathcal{H}_t := \Gamma_4(t) \in \mathbb{R}^{d \times d}$ is positive semi-definite.
- (iv) There exists a constant $b_V > 0$ such that, for a.e. $t \in [0, T]$,

$$\frac{d}{dt} c_t = -\mathcal{H}_t c_t \tag{57a}$$

$$\frac{d}{dt} V(w_t) \leq -U(w_t) \tag{57b}$$

with $V(w_t) = \frac{1}{2} \|\bar{f}(w_t)\|^2$ and $U(w_t) = b_V \|\mathcal{A}_t^\top c_t\|^2$. \square

An alert reader will notice that we have *not* obtained the desired ODE limit, since (57a) may differ from the ODE solution given in Lemma A.6 (i). In particular, we do not know if \mathcal{A}_t coincides with $A(w_t)$ (where $A(\theta)$ is defined in (20) using a particular Q^θ -greedy policy), and we do not know if \mathcal{H}_t coincides with

$$A(w_t)[\varepsilon I + A(w_t)^\top A(w_t)]^{-1} A(w_t)^\top$$

We preserve the essential drift condition (57b), which leads to a simple proof of the main result:

Proof of Thm. 2.1. Prop. A.12 (i) and (ii) justify the assertion that $U(w_t) := b_V \|\mathcal{A}_t^\top c_t\|^2$ is in fact a function of w_t . Under (A2) we see that Assumption (i) of Lemma A.10 holds, and (57b) implies Assumption (ii) of the lemma.

For given $\eta > 0$, we may choose $T \geq V(w_0)B(\eta)$, so that $V(w_T) \leq \eta$ for any sub-sequential limit. It then follows that $\limsup_{n \rightarrow \infty} V(\theta_n) \leq \eta$. Since $\eta > 0$ is arbitrary, it follows that $V(w_T) \equiv 0$; that is, $\lim_{t \rightarrow \infty} \bar{f}(\theta_n) = 0$ as claimed. \square

A.5.2 Analysis of $\{\hat{A}_n\}$ over the fast time scale

The goal in this subsection is to show that \hat{A}_n is close to the set $\mathcal{A}(\theta_n)$ with n sufficiently large. An explicit representation of $\mathcal{A}(\theta)$ is given in the following: denote, for any $\theta \in \mathbb{R}^d$ and any (possibly randomized) policy ϕ , the random $d \times d$ matrix:

$$A_{n+1}(\theta, \phi) = [\gamma \partial_\theta Q^\theta(X_{n+1}, \phi(X_{n+1})) - \partial_\theta Q^\theta(X_n, U_n)] \zeta(\theta, X_n, U_n)$$

If ϕ is Q^θ -greedy, meaning

$$Q^\theta(x, \phi(x)) = \underline{Q}^\theta(x), \quad x \in \mathbb{X},$$

then a generalized subgradient of the function f in (19b) is given by

$$A_{n+1}(\theta, \phi) + \mathcal{D}(\theta, \Phi_{n+1}) \partial_\theta \zeta(\theta, X_n, U_n)$$

Lemma A.13. *If (A1)-(A2) hold, and if ζ is non-negative, then the set $\mathcal{A}(\theta)$ defined in (26) admits the representation,*

$$\mathcal{A}(\theta) = \{E_\varpi[A_{n+1}(\theta, \tilde{\phi}_{n+1}^\theta) + \mathcal{D}(\theta, \Phi_{n+1}) \partial_\theta \zeta(\theta, X_n, U_n)] : \tilde{\phi}_{n+1}^\theta \text{ is } Q^\theta\text{-greedy}\}$$

where $\tilde{\phi}_{n+1}^\theta$ ranges over all Q^θ -greedy randomized policies. \square

A key implication of the non-negativity assumption (54) is the following:

Lemma A.14. *Under Assumptions (A1)-(A2) and (54), there exists $b_T < \infty$ such that, for all $n \geq 1$ and all vectors $v \in \mathbb{R}^d$, $\|v\| \leq 1$,*

$$f(\theta_n + v, \Phi_{n+1}) \geq f(\theta_n, \Phi_{n+1}) + A_{n+1}v - b_T \|v\|^2 \mathbf{1} \quad (58)$$

where the inequality is component-wise, A_{n+1} is defined in (21a), and $\mathbf{1}$ denotes the vector of all ones. In particular, when $Q^\theta = \psi^\top \theta$, we have $b_T = 0$:

$$f(\theta_n + v, \Phi_{n+1}) \geq f(\theta_n, \Phi_{n+1}) + A_{n+1}v \quad (59)$$

Proof. The proof is based on the Taylor series expansion. With $z := (x', x, u', u)$, define $g : \mathbb{R}^d \times \mathbb{Z} \times \mathbb{U} \rightarrow \mathbb{R}$ by

$$g(\theta, z, u^\circ) := r(x, u) + \gamma Q^\theta(x', u^\circ) - Q^\theta(x, u) \quad (60)$$

By (A2), g admits the Taylor series expansion at each $\|\theta\| \leq b_\theta$:

$$g(\theta + v, z, u^\circ) = g(\theta, z, u^\circ) + \partial_\theta g(\theta, z, u^\circ)v + O(\|v\|^2)$$

Recall that $\mathcal{D}(\theta, z) := g(\theta, z, \phi^\theta(x')) = \max_{u^\circ} g(\theta, z, u^\circ)$ and the state-input space is finite,

$$\begin{aligned} \mathcal{D}(\theta_n + v, \Phi_{n+1}) &= \max_{u^\circ} g(\theta_n + v, \Phi_{n+1}, u^\circ) \\ &= \max_{u^\circ} g(\theta_n, \Phi_{n+1}, u^\circ) + \partial_\theta g(\theta_n, \Phi_{n+1}, u^\circ)v + O(\|v\|^2) \\ &\geq \mathcal{D}(\theta_n, \Phi_{n+1}) + \partial_\theta g(\theta_n, \Phi_{n+1}, \phi^{\theta_n}(X_{n+1}))v + O(\|v\|^2) \end{aligned} \quad (61)$$

Denote $\zeta_n(\theta) := \zeta(\theta, X_n, U_n)$. Another Taylor series expansion of ζ at θ_n gives

$$\zeta_n(\theta_n + v) = \zeta_n(\theta_n) + \partial_\theta \zeta_n(\theta_n)v + O(\|v\|^2) \quad (62)$$

We next recall that $f(\theta_n, \Phi_{n+1}) = \zeta_n(\theta_n) \mathcal{D}(\theta_n, \Phi_{n+1})$,

$$\begin{aligned} f(\theta_n + v, \Phi_{n+1}) - f(\theta_n, \Phi_{n+1}) &= \zeta_n(\theta_n) \{ \mathcal{D}(\theta_n + v, \Phi_{n+1}) - \mathcal{D}(\theta_n, \Phi_{n+1}) \} \\ &\quad + \{ \zeta_n(\theta_n + v) - \zeta_n(\theta_n) \} \mathcal{D}(\theta_n, \Phi_{n+1}) \\ &\quad + \{ \zeta_n(\theta_n + v) - \zeta_n(\theta_n) \} \{ \mathcal{D}(\theta_n + v, \Phi_{n+1}) - \mathcal{D}(\theta_n, \Phi_{n+1}) \} \end{aligned}$$

By (61) and the non-negativity assumption (54),

$$\zeta_n(\theta_n) \{ \mathcal{D}(\theta_n + v, \Phi_{n+1}) - \mathcal{D}(\theta_n, \Phi_{n+1}) \} \geq \zeta_n(\theta_n) \partial_\theta g(\theta_n, \Phi_{n+1}, \phi^{\theta_n}(X_{n+1}))v + O(\|v\|^2)$$

Similarly, from (62),

$$\begin{aligned} \{\zeta_n(\theta_n + v) - \zeta_n(\theta_n)\} \mathcal{D}(\theta_n, \Phi_{n+1}) &= \{\partial_\theta \zeta_n(\theta_n) v + O(\|v\|^2)\} \mathcal{D}(\theta_n, \Phi_{n+1}) \\ &\geq \mathcal{D}(\theta_n, \Phi_{n+1}) \partial_\theta \zeta_n(\theta_n) v + O(\|v\|^2) \end{aligned}$$

By (A2) once more, both ζ and \mathcal{D} are Lipschitz continuous in θ ,

$$\|\zeta_n(\theta_n + v) - \zeta_n(\theta_n)\| \|\mathcal{D}(\theta_n + v, \Phi_{n+1}) - \mathcal{D}(\theta_n, \Phi_{n+1})\| = O(\|v\|^2)$$

Consequently,

$$\begin{aligned} f(\theta_n + v, \Phi_{n+1}) - f(\theta_n, \Phi_{n+1}) \\ \geq \{\zeta_n(\theta_n) \partial_\theta g(\theta_n, \Phi_{n+1}, \phi^{\theta_n}(X_{n+1})) + \mathcal{D}(\theta_n, \Phi_{n+1}) \partial_\theta \zeta_n(\theta_n)\} v + O(\|v\|^2) \end{aligned}$$

The proof is completed by realizing that A_{n+1} defined in (21a) can be expressed

$$A_{n+1} = \zeta_n(\theta_n) \partial_\theta g(\theta_n, \Phi_{n+1}, \phi^{\theta_n}(X_{n+1})) + \mathcal{D}(\theta_n, \Phi_{n+1}) \partial_\theta \zeta_n(\theta_n)$$

□

Define the *fast time scale*, over which the matrix gain sequence $\{\hat{A}_n\}$ is updated,

$$t_n = \sum_{i=1}^n \beta_i = \sum_{i=1}^n \frac{1}{i^\rho}, \quad n \geq 1, \quad t_0 = 0, \quad \rho \in (0.5, 1) \quad (63)$$

Define the time process $\{\bar{\mathcal{A}}_t : t \geq 0\}$ with $\bar{\mathcal{A}}_{t_n} = \hat{A}_n$ for those values t_n , with the definition extended to \mathbb{R}_+ via linear interpolation. Note that this definition of $\{\bar{\mathcal{A}}_t : t \geq 0\}$ is used only in this subsection to analyze $\{\hat{A}_n\}$. For each $n \geq 1$, define the associated time block: $[t_{m(n)}, t_n]$ where $m(n) = \min\{j : t_j + \ln(n) \geq t_n\}$. Some properties of this fast time scale setting are collected in the following:

Lemma A.15. *The follow hold:*

- (i) $\ln(n) - 1 < t_n - t_{m(n)} \leq \ln(n)$.
- (ii) *There exists $N_s \geq 1$ such that for $n \geq N_s$, $m(n) + 1 \geq \rho^{1/(1-\rho)}(n+1)$.*
- (iii) $\lim_{n \rightarrow \infty} \max_{m(n) \leq k \leq n} \|\theta_k - \theta_n\| = 0$.

Proof. (i) follows directly from the definition.

By (63),

$$\begin{aligned} t_n - t_{m(n)} &= \sum_{i=m(n)+1}^n \frac{1}{i^\rho} \geq \int_{m(n)+1}^{n+1} \frac{1}{\tau^\rho} d\tau \\ &= (1-\rho)^{-1} [(n+1)^{1-\rho} - (m(n)+1)^{1-\rho}] \end{aligned} \quad (64)$$

Since $\ln(n) \geq t_n - t_{m(n)}$, we have

$$(1-\rho) \ln(n) \geq (n+1)^{1-\rho} - (m(n)+1)^{1-\rho}$$

There exists $N_s \geq 1$ such that $(n+1)^{1-\rho} \geq \ln(n)$ for $n \geq N_s$. Hence,

$$(1-\rho)(n+1)^{1-\rho} \geq (n+1)^{1-\rho} - (m(n)+1)^{1-\rho}, \quad n \geq N_s$$

which proves (ii).

By (64),

$$\begin{aligned} (1-\rho)^{-1} [(n+1)^{1-\rho} - (k+1)^{1-\rho}] &\leq (1-\rho)^{-1} [(n+1)^{1-\rho} - (m(n)+1)^{1-\rho}] \\ &\leq \ln(n) \end{aligned}$$

Multiplying each side of above inequality by $(1-\rho)(k+1)^{\rho-1}$ gives

$$\left(\frac{n+1}{k+1}\right)^{1-\rho} - 1 \leq (1-\rho)(k+1)^{\rho-1} \ln(n) \leq (1-\rho)(m(n)+1)^{\rho-1} \ln(n)$$

By the inequality $\ln(1+x) \leq x$ for $x > -1$,

$$(1-\rho) \ln\left(\frac{n+1}{k+1}\right) \leq \ln(1+(1-\rho)(m(n)+1)^{\rho-1} \ln(n)) \leq (1-\rho)(m(n)+1)^{\rho-1} \ln(n)$$

Given $m(n)+1 \geq \rho^{1/(1-\rho)}(n+1)$ in (ii),

$$\ln\left(\frac{n+1}{k+1}\right) \leq \rho^{-1} \ln(n)(n+1)^{\rho-1} \quad (65)$$

The parameter vector θ_n updated by (21d) can be expressed

$$\theta_n = \theta_k + \sum_{i=k+1}^n \alpha_i G_i f(\theta_{i-1}, \Phi_i), \quad m(n) \leq k < n$$

We can find a constant $b_f < \infty$ such that $\sup_n \|G_{n+1} f(\theta_n, \Phi_{n+1})\| \leq b_f$ for almost every $\omega \in \Omega$. With $\alpha_i \equiv 1/i$,

$$\|\theta_n - \theta_k\| \leq b_f \sum_{i=k+1}^n \alpha_i \leq b_f \int_k^n \frac{1}{\tau} d\tau \leq b_f \ln\left(\frac{n}{k}\right)$$

By (65), for $n \geq N_s$,

$$\begin{aligned} \|\theta_n - \theta_k\| &\leq b_f \left| \ln\left(\frac{n}{k}\right) - \ln\left(\frac{n+1}{k+1}\right) \right| + b_f \rho^{-1} \ln(n)(n+1)^{\rho-1} \\ &\leq b_f \left| \ln\left(1 - \frac{1}{n+1}\right) + \ln\left(1 + \frac{1}{k}\right) \right| + b_f \rho^{-1} \ln(n)(n+1)^{\rho-1} \\ &\leq b_f \frac{1}{k} + b_f \rho^{-1} \ln(n)(n+1)^{\rho-1} \\ &\leq b_f \frac{1}{\rho^{1/(1-\rho)}(n+1) - 1} + b_f \rho^{-1} \ln(n)(n+1)^{\rho-1} \end{aligned} \quad (66)$$

where the last inequality holds given $k \geq m(n) \geq \rho^{1/(1-\rho)}(n+1) - 1$. Therefore, $\max_{m(n) \leq k \leq n} \|\theta_k - \theta_n\| \rightarrow 0$ as $n \rightarrow \infty$. \square

Proposition A.16. *Under Assumptions (A1)-(A2) and (54), the following hold for all $v \in \mathbb{R}^d$, $\|v\| \leq 1$, and all $k \in \mathbb{Z}$ between $m(n)$ and n :*

(i)

$$\sum_{i=k+1}^n \beta_i [f(\theta_{i-1} + v, \Phi_i) - f(\theta_{i-1}, \Phi_i) + b_T \|v\|^2 \mathbf{1}] \geq \hat{A}_n v - \hat{A}_k v + \sum_{i=k+1}^n \beta_i \hat{A}_{i-1} v \quad (67)$$

(ii) For any $t \in [t_{m(n)}, t_n]$,

$$\bar{\mathcal{A}}_{t_n} v - \bar{\mathcal{A}}_t v + \int_t^{t_n} \bar{\mathcal{A}}_\tau v d\tau \leq (t_n - t) [\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1}] + o(1), \quad n \rightarrow \infty \quad (68)$$

where $o(1) \rightarrow 0$ as $n \rightarrow \infty$, uniformly in v .

Proof. By (58), for each $n \geq 1$,

$$f(\theta_n + v, \Phi_{n+1}) \geq f(\theta_n, \Phi_{n+1}) + A_{n+1} v - b_T \|v\|^2 \mathbf{1}, \quad v \in \mathbb{R}^d$$

Consequently,

$$\sum_{i=k+1}^n \beta_i [f(\theta_{i-1} + v, \Phi_i) - f(\theta_{i-1}, \Phi_i) + b_T \|v\|^2 \mathbf{1}] \geq \sum_{i=k+1}^n \beta_i A_i v, \quad m(n) \leq k \leq n$$

The gain matrix \hat{A}_n updated by (21b) can be expressed

$$\hat{A}_n = \hat{A}_k + \sum_{i=k+1}^n \beta_i A_i - \sum_{i=k+1}^n \beta_i \hat{A}_{i-1}, \quad m(n) \leq k \leq n$$

Therefore, $\sum_{i=k+1}^n \beta_i A_i v = \hat{A}_n v - \hat{A}_k v + \sum_{i=k+1}^n \beta_i \hat{A}_{i-1} v$. This proves (i).

Now consider the sum $\sum_{i=k+1}^n \beta_i f(\theta_{i-1} + v, \Phi_i)$ with $m(n) \leq k \leq n$. We first rewrite it in the suggestive form

$$\sum_{i=k+1}^n \beta_i f(\theta_{i-1} + v, \Phi_i) = \sum_{i=k+1}^n \beta_i [f(\theta_{i-1} + v, \Phi_i) - f(\theta_n + v, \Phi_i)] + \sum_{i=k+1}^n \beta_i f(\theta_n + v, \Phi_i)$$

By the Lipschitz continuity of f in θ and Lemma A.15 (iii), the first sum on the right hand side goes to 0 uniformly in k as $n \rightarrow \infty$. The second sum can be expressed

$$\sum_{i=k+1}^n \beta_i f(\theta_n + v, \Phi_i) = (t_n - t_k) \bar{f}(\theta_n + v) + \sum_{i=k+1}^n \beta_i [f(\theta_n + v, \Phi_i) - \bar{f}(\theta_n + v)]$$

Each term in the sum on the right side has zero-mean under the stationary pmf of (\mathbf{X}, \mathbf{U}) . It goes to zero *a.s.* for $m(n) \leq k \leq n$ as $n \rightarrow \infty$ [5, Part II, Section 1.4.6, Proposition 7]. We then obtain

$$\max_{m(n) \leq k \leq n} \left\| (t_n - t_k) \bar{f}(\theta_n + v) - \sum_{i=k+1}^n \beta_i f(\theta_{i-1} + v, \Phi_i) \right\| = o(1), \quad n \rightarrow \infty \quad (69)$$

Since the process $\{\bar{\mathcal{A}}_t : t \geq 0\}$ is linearly interpolated between discrete values,

$$\int_{t_k}^{t_n} \bar{\mathcal{A}}_\tau v d\tau = \frac{1}{2} \sum_{i=k+1}^n \beta_i [\hat{A}_i + \hat{A}_{i-1}] v = \sum_{i=k+1}^n \beta_i \hat{A}_{i-1} v + \frac{1}{2} \sum_{i=k+1}^n \beta_i [\hat{A}_i - \hat{A}_{i-1}] v$$

where the second sum on the right hand side can be rewritten as

$$\sum_{i=k+1}^n \beta_i [\hat{A}_i - \hat{A}_{i-1}] v = -\beta_{k+1} \hat{A}_k v + \beta_n \hat{A}_{n+1} v + \sum_{i=k+1}^{n-1} [\beta_i - \beta_{i+1}] \hat{A}_i v$$

which goes to zero as $n \rightarrow \infty$ given $\sup_n \|\hat{A}_n\| < \infty$ and $\beta_i - \beta_{i+1} \approx \rho i^{-1} \beta_i$. Therefore,

$$\max_{m(n) \leq k \leq n} \left\| \sum_{i=k+1}^n \beta_i \hat{A}_{i-1} v - \int_{t_k}^{t_n} \bar{\mathcal{A}}_\tau v d\tau \right\| = o(1), \quad n \rightarrow \infty \quad (70)$$

Combining (i) with (69) and (70) gives, for $t \in \{t_k : m(n) \leq k \leq n\}$,

$$\bar{\mathcal{A}}_{t_n} v - \bar{\mathcal{A}}_t v + \int_t^{t_n} \bar{\mathcal{A}}_\tau v d\tau \leq (t_n - t) [\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1}] + o(1) \quad (71)$$

For any $t \in [t_{m(n)}, t_n]$, denote $k = \max\{j : t_j \leq t\}$. Letting $\delta = (t - t_k)/(t_{k+1} - t_k)$, we have

$$\bar{\mathcal{A}}_t v = (1 - \delta) \bar{\mathcal{A}}_{t_k} v + \delta \bar{\mathcal{A}}_{t_{k+1}} v$$

Then,

$$\begin{aligned} (1 - \delta) \{ \bar{\mathcal{A}}_{t_n} v - \bar{\mathcal{A}}_{t_k} v + \int_{t_k}^{t_n} \bar{\mathcal{A}}_\tau v d\tau \} &\leq (1 - \delta) (t_n - t_k) [\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1}] + o(1) \\ \delta \{ \bar{\mathcal{A}}_{t_n} v - \bar{\mathcal{A}}_{t_{k+1}} v + \int_{t_{k+1}}^{t_n} \bar{\mathcal{A}}_\tau v d\tau \} &\leq \delta (t_n - t_{k+1}) [\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1}] + o(1) \end{aligned}$$

Combining above two inequalities gives

$$\bar{\mathcal{A}}_{t_n} v - \bar{\mathcal{A}}_t v + \int_t^{t_n} \bar{\mathcal{A}}_\tau v d\tau \leq (t_n - t) [\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1}] + o(1)$$

□

Recall the constant $b_T > 0$ introduced in Lemma A.14. For fixed matrix $\hat{A} \in \mathbb{R}^{d \times d}$ and vector $\theta \in \mathbb{R}^d$, define the function $\text{dist}_{\mathcal{N}} : \mathbb{R}^{d \times d} \times \mathbb{R}^d \rightarrow \mathbb{R}$ by

$$\text{dist}_{\mathcal{N}}(\hat{A}, \theta) = \sup_{\|v\| \leq 1} \left\{ \max_i [\hat{A}v - (\bar{f}(\theta + v) - \bar{f}(\theta))]_i - b_T \|v\|^2 \right\} \quad (72)$$

This measures how well $\hat{A}v$ approximates the directional derivative $f'(\theta; v)$ for v in the unit ball. It is non-negative since $v = 0$ is feasible in the supremum in (72). It is also continuous in both arguments:

Proposition A.17. *Under Assumptions (A1)-(A2) and (54), the function $\text{dist}_{\mathcal{N}}$ defined in (72) satisfies:*

- (i) *For fixed \hat{A} and θ , the supremum in (72) is achieved.*
- (ii) *$\text{dist}_{\mathcal{N}}(\hat{A}, \theta)$ is non-negative and Lipschitz continuous in both \hat{A} and θ .*
- (iii) *If $\text{dist}_{\mathcal{N}}(\hat{A}, \theta) = 0$, then the following hold: $\hat{A} \in \mathcal{A}(\theta)$, and*

$$\text{If } \bar{f}'(\theta; v) = -\bar{f}'(\theta; -v) \text{ for some } \|v\| \leq 1, \text{ then } \hat{A}v = \bar{f}'(\theta; v).$$

Proof. With fixed \hat{A} and θ , $\max_i [\hat{A}v - (\bar{f}(\theta + v) - \bar{f}(\theta))]_i$ is Lipschitz continuous with respect to v by Lemma A.8 (i). Since the set $\{v : \|v\| \leq 1\}$ is compact, the supremum is achieved.

For (ii), consider $\hat{A} \neq \hat{A}'$, while θ is fixed. Let v^*, i^* maximize $[\hat{A}v - (\bar{f}(\theta + v) - \bar{f}(\theta))]_i - b_T \|v\|^2$. We have

$$\begin{aligned} \text{dist}_{\mathcal{N}}(\hat{A}, \theta) - \text{dist}_{\mathcal{N}}(\hat{A}', \theta) &\leq [\hat{A}v^* - (\bar{f}(\theta + v^*) - \bar{f}(\theta))]_{i^*} - [\hat{A}'v^* - (\bar{f}(\theta + v^*) - \bar{f}(\theta))]_{i^*} \\ &\leq \|\hat{A} - \hat{A}'\|_1 \|v^*\|_1 \end{aligned}$$

Therefore, $\text{dist}_{\mathcal{N}}(\hat{A}, \theta)$ is Lipschitz continuous in \hat{A} . The same argument implies the Lipschitz continuity of $\text{dist}_{\mathcal{N}}(\hat{A}, \theta)$ in θ .

For (iii), the first claim follows from the definition of $\mathcal{A}(\theta)$ in (26). By the definition of directional derivative,

$$\bar{f}'(\theta; v) = \bar{f}(\theta + v) - \bar{f}(\theta) + o(\|v\|) \quad (73)$$

where $o(s)/s \rightarrow 0$ as $s \downarrow 0$. Given $\text{dist}_{\mathcal{N}}(\hat{A}, \theta) = 0$, we have for each $v \in \mathbb{R}^d$,

$$\begin{aligned} \hat{A}v &\leq \bar{f}(\theta + v) - \bar{f}(\theta) + b_T \|v\|^2 \mathbf{1} = \bar{f}'(\theta; v) + o(\|v\|) \\ -\hat{A}v &\leq \bar{f}(\theta - v) - \bar{f}(\theta) + b_T \|v\|^2 \mathbf{1} = \bar{f}'(\theta; -v) + o(\|v\|) \end{aligned}$$

Using $\bar{f}'(\theta; -v) = -\bar{f}'(\theta; v)$ gives

$$\bar{f}'(\theta; v) - o(\|v\|) \leq \hat{A}v \leq \bar{f}'(\theta; v) + o(\|v\|)$$

With $\bar{f}'(\theta; sv)/s = \bar{f}'(\theta; v)$ for $s > 0$, replace v by sv in the above inequality and divide:

$$\bar{f}'(\theta; v) - \frac{o(s\|v\|)}{s} \leq \hat{A}v \leq \bar{f}'(\theta; v) + \frac{o(s\|v\|)}{s}$$

Letting $s \downarrow 0$ gives $\hat{A}v = \bar{f}'(\theta; v)$. □

Proposition A.18. *Under Assumptions (A1)-(A2) and (54),*

- (i) *The component-wise inequality holds:*

$$\hat{A}_n v \leq \bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1} + o(1), \quad n \rightarrow \infty \quad (74)$$

where $o(1) \rightarrow 0$ as $n \rightarrow \infty$ uniformly in $\|v\| \leq 1$.

- (ii) $\lim_{n \rightarrow \infty} \text{dist}_{\mathcal{N}}(\hat{A}_n, \theta_n) = 0$ a.s..

(iii) Let $\{\theta_{n_k}\}$ be a subsequence of $\{\theta_n\}$ that converges to some $\theta^\circ \in \mathbb{R}^d$ a.s.. Then,

$$\lim_{k \rightarrow \infty} \text{dist}(\hat{A}_{n_k}, \mathcal{A}(\theta^\circ)) = 0, \quad a.s. \quad (75)$$

where $\text{dist}(\hat{A}_{n_k}, \mathcal{A}(\theta^\circ))$ denotes the Euclidean distance between \hat{A}_{n_k} and the set $\mathcal{A}(\theta^\circ)$.

Proof. For fixed n and $v \in \mathbb{R}^d$, let $\mathcal{U} : [t_{m(n)}, t_n] \rightarrow \mathbb{R}^d$ denote the solution of the following linear integral equation

$$\mathcal{U}_t = \mathcal{U}_{t_{m(n)}} - \int_{t_{m(n)}}^t \mathcal{U}_\tau d\tau + (t - t_{m(n)})[\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2], \quad \mathcal{U}_{t_{m(n)}} = \bar{\mathcal{A}}_{t_{m(n)}} v \quad (76)$$

With n fixed, $\delta_n \triangleq \max_i |o(1)_i|$ in (68) can be viewed as a positive constant. We claim that $\bar{\mathcal{A}}_{t_n} v \leq \mathcal{U}_{t_n} + \mathbf{1}\delta_n$. Suppose the claim is not true. Then $[\bar{\mathcal{A}}_{t_n} v]_i > [\mathcal{U}_{t_n}]_i + \delta_n$ for some index i between 1 and d . Because $\bar{\mathcal{A}}_t v$ and \mathcal{U}_t are both continuous functions over $[t_{m(n)}, t_n]$ and $\bar{\mathcal{A}}_{t_{m(n)}} v = \mathcal{U}_{t_{m(n)}}$, there exists $t \in [t_{m(n)}, t_n]$ such that $[\bar{\mathcal{A}}_t v]_i = [\mathcal{U}_t]_i$ and $[\bar{\mathcal{A}}_\tau v]_i > [\mathcal{U}_\tau]_i$ for $\tau \in (t, t_n)$. Consequently, combining (68) and (76) gives

$$\delta_n < [\bar{\mathcal{A}}_{t_n} v - \mathcal{U}_{t_n}]_i \leq [\bar{\mathcal{A}}_t v - \mathcal{U}_t]_i - \int_t^{t_n} [\bar{\mathcal{A}}_\tau v - \mathcal{U}_\tau]_i d\tau + \delta_n < \delta_n$$

which is a contradiction. Therefore, $\bar{\mathcal{A}}_{t_n} v \leq \mathcal{U}_{t_n} + \mathbf{1}\delta_n$.

The integral equation (76) has the solution,

$$\mathcal{U}_t = \exp(t_{m(n)} - t) \mathcal{U}_{t_{m(n)}} + (1 - \exp(t_{m(n)} - t))[\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2], \quad t \in [t_{m(n)}, t_n]$$

Consequently,

$$\begin{aligned} \hat{A}_n v &\leq \bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1} + \delta_n \mathbf{1} \\ &\quad + \exp(t_{m(n)} - t_n) [\mathcal{U}_{t_{m(n)}} - (\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2)] \end{aligned}$$

By Lemma A.15 (i), we have $t_{m(n)} - t_n < -\ln(n) + 1$ and hence $\exp(t_{m(n)} - t_n) < e/n$. Therefore,

$$\begin{aligned} \|\exp(t_{m(n)} - t_n) [\mathcal{U}_{t_{m(n)}} - (\bar{f}(\theta_n + v) - \bar{f}(\theta_n) + b_T \|v\|^2 \mathbf{1})]\| &\leq \frac{e}{n} [b_A + b_L + b_T] \|v\| \\ &\leq \frac{e}{n} [b_A + b_L + b_T] \end{aligned}$$

which goes to zero as $n \rightarrow \infty$. This proves (i), and (ii) follows by the definition of $\text{dist}_{\mathcal{N}}$.

We prove (iii) by contradiction: Suppose (75) does not hold. Then there exists a constant $\delta > 0$ and a subsequence $\{\hat{A}_{n_k}\}$ such that $\text{dist}(\hat{A}_{n_k}, \mathcal{A}(\theta^\circ)) \geq \delta$ for each k . Without loss of generality, the subsequence is convergent, with limit \hat{A}° satisfying $\text{dist}(\hat{A}^\circ, \mathcal{A}(\theta^\circ)) \geq \delta$. However, combining statement (i) and Prop. A.17 (iii) gives

$$\text{dist}(\hat{A}^\circ, \mathcal{A}(\theta^\circ)) = 0$$

which is a contradiction. \square

A.5.3 Proofs of Prop. A.11 and Prop. A.12

In this subsection, the time processes involved all refer to those defined in Section A.5.1 with respect to the slow time scale (55).

Proof of Prop. A.11. The Lipschitz continuity of $\{\bar{w}_t\}$ and $\{\bar{c}_t\}$ follows directly from boundedness of $\{\theta_n\}$.

At a point of differentiability, let $v_t = \frac{d}{dt} \bar{w}_t = \mathcal{G}_t f(\theta_{[t]}, \Phi_{[t]+1})$ and recall that $\sup_t \|v_t\| \leq b_f$. Whenever exists, the derivative of \bar{c}_t may be represented as the directional derivative of $\bar{f}(\bar{w}_t)$ along direction v_t :

$$\begin{aligned} \frac{d}{dt} \bar{c}_t &= \lim_{s \rightarrow 0} \frac{\bar{f}(\bar{w}_{t+s}) - \bar{f}(\bar{w}_t)}{s} = \lim_{s \downarrow 0} \frac{\bar{f}(\bar{w}_{t+s}) - \bar{f}(\bar{w}_t)}{s} = \bar{f}'(\bar{w}_t; v_t) \\ &= \lim_{s \uparrow 0} \frac{\bar{f}(\bar{w}_{t+s}) - \bar{f}(\bar{w}_t)}{s} = -\bar{f}'(\bar{w}_t; -v_t) \end{aligned} \quad (77)$$

Prop. A.17 (ii) combined with Prop. A.18 (ii) gives

$$\lim_{t \rightarrow \infty} \text{dist}_{\mathcal{N}}(\bar{\mathcal{A}}_t, \bar{w}_t) \leq 0, \quad a.s. \quad (78)$$

Let $\eta_t := \max(1/t, \text{dist}_{\mathcal{N}}(\bar{\mathcal{A}}_t, \bar{w}_t))$, satisfying $\eta_t > 0$ and $\eta_t \rightarrow 0$ as $t \rightarrow \infty$. There exists $T_\bullet < \infty$ a.s. such for $t \geq T_\bullet$,

$$\begin{aligned} \bar{\mathcal{A}}_t v_t - \bar{f}'(\bar{w}; v_t) &= \frac{1}{\sqrt{\eta_t}} [\bar{\mathcal{A}}_t \sqrt{\eta_t} v_t - \bar{f}'(\bar{w}_t; \sqrt{\eta_t} v_t)] \\ &= \frac{1}{\sqrt{\eta_t}} [\bar{\mathcal{A}}_t \sqrt{\eta_t} v_t - [\bar{f}(\bar{w}_t + \sqrt{\eta_t} v_t) - \bar{f}(\bar{w}_t)]] + o(\|v_t\|) \\ &\leq (1 + b_T b_f) \sqrt{\eta_t} \mathbf{1} + o(\|v_t\|) \end{aligned} \quad (79)$$

where the second equality follows from (77) and the last inequality holds given $\text{dist}_{\mathcal{N}}(\bar{\mathcal{A}}_t, \bar{w}_t) \leq \eta_t$ and $\|v_t\|$ is uniformly bounded by b_f .

At points of differentiability, we apply $\bar{f}'(\bar{w}_t; v_t) = -\bar{f}'(\bar{w}_t; -v_t)$ from (77):

$$-\bar{\mathcal{A}}_t v_t + \bar{f}'(\bar{w}; v_t) \leq (1 + b_T b_f) \sqrt{\eta_t} \mathbf{1} + o(\|v_t\|)$$

Consequently,

$$\|\bar{\mathcal{A}}_t \frac{d}{dt} \bar{w}_t - \frac{d}{dt} \bar{c}_t\|_\infty \leq (1 + b_T b_f) \sqrt{\eta_t} + o(\|v_t\|)$$

where $\|\cdot\|_\infty$ denotes the infinity norm. The right hand side of above inequality is bounded and converges to zero as $t \rightarrow \infty$. Since the derivatives of \bar{w}_t and \bar{c}_t exist a.e., we have for each $T > 0$,

$$\int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t \frac{d}{dt} \bar{w}_t - \frac{d}{dt} \bar{c}_t\|_\infty dt \leq \int_{T_0}^{T_0+T} (1 + b_T b_f) \sqrt{\eta_t} + o(\|v_t\|) dt$$

The desired result follows from Dominated Convergence Theorem.

Part (ii) is obtained from (i):

$$\begin{aligned} \bar{c}_{T_0+t} &= \bar{c}_{T_0} + \int_{T_0}^{T_0+t} \frac{d}{d\tau} \bar{c}_\tau d\tau \\ &= \bar{c}_{T_0} + \int_{T_0}^{T_0+t} \bar{\mathcal{A}}_\tau \bar{\mathcal{G}}_\tau f(\theta_{[\tau]}, \Phi_{[\tau]+1}) d\tau + o(1), \quad T_0 \rightarrow \infty \\ &= \bar{c}_{T_0} + \int_{T_0}^{T_0+t} \bar{\mathcal{A}}_\tau \bar{\mathcal{G}}_\tau \bar{f}(\bar{w}_\tau) d\tau + o(1), \quad T_0 \rightarrow \infty \end{aligned}$$

where the last equality follows from standard ODE arguments for stochastic approximation [5].

For (iii), $\|\bar{c}_t\|^2$ is Lipschitz continuous in t given boundedness of $\{\theta_n\}$. Hence by the same argument in (ii),

$$\begin{aligned} \|\bar{c}_{T_0+t}\|^2 &= \|\bar{c}_{T_0}\|^2 + 2 \int_{T_0}^{T_0+t} \bar{c}_\tau^\top \frac{d}{d\tau} \bar{c}_\tau d\tau \\ &= \|\bar{c}_{T_0}\|^2 + 2 \int_{T_0}^{T_0+t} \bar{c}_\tau^\top \bar{\mathcal{A}}_\tau \bar{\mathcal{G}}_\tau \bar{c}_\tau d\tau + o(1), \quad T_0 \rightarrow \infty \end{aligned}$$

□

Proof of Prop. A.12. (i) follows from the Lipschitz continuity of \bar{f} .

Let $\{T_k\}$ be a sequence such that $\bar{\Gamma}^{T_k} \rightarrow \Gamma$ for each of the four components: $\bar{\Gamma}_i^{T_k}$, $1 \leq i \leq 4$. Since $\bar{\Gamma}_3^{T_k} \rightarrow \mathcal{A}$ weakly in $L_2([0, T]; \mathbb{R}^{d \times d})$ as $k \rightarrow \infty$, by the Banach-Saks theorem, there exists a subsequence $\{T_{n_k}\}$ such that

$$\frac{1}{N} \sum_{k=1}^N \bar{\Gamma}_3^{T_{n_k}}(t) \rightarrow \mathcal{A}_t, \quad a.e. t \in [0, T], \quad N \rightarrow \infty$$

Without loss of generality, we can modifying \mathcal{A}_t on a Lebesgue-null set such that the convergence above is pointwise. We also have $\bar{\Gamma}_1^{T_{n_k}}(t) \rightarrow w_t$ as $k \rightarrow \infty$ for each $t \in [0, T]$. By Prop. A.18 (ii),

$$\lim_{k \rightarrow \infty} \text{dist}(\bar{\Gamma}_3^{T_{n_k}}(t), \mathcal{A}(w_t)) = 0, \quad t \in [0, T]$$

It follows from definition (26) that the set $\mathcal{A}(\theta)$ is convex for each θ . Then,

$$\lim_{N \rightarrow \infty} \text{dist}\left(\frac{1}{N} \sum_{k=1}^N \bar{\Gamma}_3^{T_{n_k}}(t), \mathcal{A}(w_t)\right) = 0, \quad t \in [0, T]$$

Therefore, $\mathcal{A}_t \in \mathcal{A}(w_t)$ for each $t \in [0, T]$. This proves (ii).

Given that $\bar{\Gamma}_4^{T_0}$ is positive semi-definite pointwise and uniformly bounded, the same arguments establish (iii).

Since $\Gamma_2^{T_k} \rightarrow c$ uniformly over $[0, T]$ and $\Gamma_4^{T_k} \rightarrow \mathcal{H}$ weakly, $\bar{\Gamma}_4^{T_k} \bar{\Gamma}_2^{T_k}$ converges to $\mathcal{H}c : [0, T] \rightarrow \mathbb{R}^d$ weakly. The ODE (57a) follows from Prop. A.11 (ii). For (57b), since $b_\lambda := \sup_n \lambda_{\max}(\hat{A}_n^\top \hat{A}_n)$ is finite,

$$-[\varepsilon I + \hat{A}_n^\top \hat{A}_n]^{-1} \leq -\frac{1}{\varepsilon + b_\lambda} I, \quad n \geq 1$$

Combining this inequality with Prop. A.11 (iii) implies

$$\|\bar{c}_{T_0+t}\|^2 \leq \|\bar{c}_{T_0}\|^2 - \frac{2}{\varepsilon + b_\lambda} \int_{T_0}^{T_0+t} \|\bar{\mathcal{A}}_\tau^\top \bar{c}_\tau\|^2 d\tau + \mathcal{O}(1), \quad T_0 \rightarrow 0 \quad (80)$$

We can show that $\{(\bar{\Gamma}_3^{T_k})^\top \bar{\Gamma}_2^{T_k}\}$ converges weakly to $\mathcal{A}^\top c$ in $L_2([0, T]; \mathbb{R}^d)$ by the same arguments that we used to establish $\bar{\Gamma}_4^{T_k} \bar{\Gamma}_2^{T_k} \rightarrow \mathcal{H}c$ weakly. Applying [17, Theorem 2.2.1], we obtain for each $t \in [0, T]$,

$$\int_0^t \|\mathcal{A}_\tau^\top c_\tau\|^2 d\tau \leq \liminf_{k \rightarrow \infty} \int_0^t \|[\bar{\Gamma}_3^{T_k}(\tau)]^\top \bar{\Gamma}_2^{T_k}(\tau)\|^2 d\tau$$

Consequently,

$$\|c_t\|^2 \leq \|c_0\|^2 - \frac{2}{\varepsilon + b_\lambda} \int_0^t \|\mathcal{A}_\tau^\top c_\tau\|^2 d\tau$$

□

A.5.4 General eligibility vector ζ

We finally come to the general model in which (54) is relaxed. For the sake of analysis, the two functions \mathcal{D}, ζ in (19b) are assumed to be parameterized by separate parameters $\theta, \xi \in \mathbb{R}^d$: $\mathcal{D}(\theta, z), \zeta(\xi, x, u)$. This is only for clarifying calculations – in the end we do impose $\theta = \xi$. Decompose the function $\zeta : \mathbb{R}^d \times \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}^d$ into its positive and negative components: $\zeta = \zeta^+ - \zeta^-$, with $\zeta^+ = \max(\zeta, 0)$ and $\zeta^- = \max(-\zeta, 0)$. Define functions $f^+, f^- : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{Z} \rightarrow \mathbb{R}^d$ by

$$f^+(\xi, \theta, z) = \zeta^+(\xi, x, u) \mathcal{D}(\theta, z), \quad f^-(\xi, \theta, z) = \zeta^-(\xi, x, u) \mathcal{D}(\theta, z)$$

Next define functions $\bar{f}^+, \bar{f}^- : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ by

$$\bar{f}^+(\xi, \theta) = \mathbb{E}_\omega[f^+(\xi, \theta, \Phi_{n+1})], \quad \bar{f}^-(\xi, \theta) = \mathbb{E}_\omega[f^-(\xi, \theta, \Phi_{n+1})]$$

Let $\mathcal{A}^+(\theta), \mathcal{A}^-(\theta)$ denote the sets of generalized subgradients of \bar{f}^+, \bar{f}^- with respect to θ based on (26). Explicit representations of $\mathcal{A}^+(\theta)$ and $\mathcal{A}^-(\theta)$ can be obtained as in Lemma A.13. With general eligibility vector ζ , let $\mathcal{A}(\theta)$ denote the set

$$\mathcal{A}(\theta) := \{A^+ - A^- + \mathbb{E}_\omega[\mathcal{D}(\theta, \Phi_{n+1}) \partial_\xi \zeta_n(\theta)] : A^+ \in \mathcal{A}^+(\theta), A^- \in \mathcal{A}^-(\theta)\} \quad (81)$$

At each $\theta \in \mathbb{R}^d$, denote

$$\bar{f}^+(\theta; v) := \lim_{s \downarrow 0} \frac{\bar{f}^+(\theta, \theta + sv) - \bar{f}^+(\theta, \theta)}{s}$$

with $\bar{f}^-(\theta; v)$ is defined similarly. Then the directional derivative $\bar{f}'(\theta; v)$ can be expressed

$$\bar{f}'(\theta; v) = \lim_{s \downarrow 0} \frac{\bar{f}(\theta + sv) - \bar{f}(\theta)}{s} = \bar{f}^+(\theta; v) - \bar{f}^-(\theta; v) + \mathbb{E}_\varpi[\mathcal{D}(\theta, \Phi_{n+1})\partial_\xi \zeta_n]v, \quad \theta, v \in \mathbb{R}^d \quad (82)$$

Decompose A_{n+1} in (21a) as $A_{n+1} = A_{n+1}^+ - A_{n+1}^- + A_{n+1}^\zeta$:

$$\begin{aligned} A_{n+1}^+ &= \zeta_n^+ [\gamma \partial_\theta Q^\theta(X_{n+1}, \phi^{\theta_n}(X_{n+1})) - \partial_\theta Q^\theta(X_n, U_n)] \\ A_{n+1}^- &= \zeta_n^- [\gamma \partial_\theta Q^\theta(X_{n+1}, \phi^{\theta_n}(X_{n+1})) - \partial_\theta Q^\theta(X_n, U_n)] \\ A_{n+1}^\zeta &= \mathcal{D}(\theta_n, \Phi_{n+1})\partial_\xi \zeta_n \end{aligned}$$

Accordingly, the matrix gain is decomposed: $\hat{A}_{n+1} = \hat{A}_{n+1}^+ - \hat{A}_{n+1}^- + \hat{A}_{n+1}^\zeta$, and each component can be expressed in the recursive form:

$$\begin{aligned} \hat{A}_{n+1}^+ &= \hat{A}_n^+ + \beta_{n+1}[A_{n+1}^+ - \hat{A}_n^+] \\ \hat{A}_{n+1}^- &= \hat{A}_n^- + \beta_{n+1}[A_{n+1}^- - \hat{A}_n^-] \\ \hat{A}_{n+1}^\zeta &= \hat{A}_n^\zeta + \beta_{n+1}[A_{n+1}^\zeta - \hat{A}_n^\zeta] \end{aligned}$$

Analysis of $\{\hat{A}_n^+, \hat{A}_n^-, \hat{A}_n^\zeta\}$ over the fast time scale: Consider the fast time scale defined by (63). The conclusions in Section A.5.2 hold for each of $\{\hat{A}_n^+\}$ and $\{\hat{A}_n^-\}$. While $\{\hat{A}_n^\zeta\}$ can be treated using standard SA arguments since A_{n+1}^ζ is Lipschitz continuous with respect to θ_n under (A2). We obtain an extension of Prop. A.18:

Proposition A.19. *The following hold:*

(i) As $n \rightarrow \infty$,

$$\begin{aligned} \hat{A}_n^+ v &\leq \bar{f}^+(\theta_n, \theta_n + v) - \bar{f}^+(\theta_n, \theta_n) + b_T \|v\|^2 \mathbf{1} + o(1) \\ \hat{A}_n^- v &\leq \bar{f}^-(\theta_n, \theta_n + v) - \bar{f}^-(\theta_n, \theta_n) + b_T \|v\|^2 \mathbf{1} + o(1) \end{aligned}$$

where $o(1) \rightarrow 0$ as $n \rightarrow \infty$, uniformly in $\|v\| \leq 1$.

(ii) Let $\{\theta_{n_k}\}$ be a subsequence of $\{\theta_n\}$ that converges to some $\theta^\circ \in \mathbb{R}^d$ a.s.. Then,

$$\lim_{k \rightarrow \infty} \text{dist}(\hat{A}_{n_k}^+, \mathcal{A}^+(\theta^\circ)) = 0, \quad \lim_{k \rightarrow \infty} \text{dist}(\hat{A}_{n_k}^-, \mathcal{A}^-(\theta^\circ)) = 0, \quad a.s.$$

(iii) $\hat{A}_n^\zeta = \mathbb{E}_\varpi[\mathcal{D}(\theta_n, \Phi_{n+1})\partial_\xi \zeta_n] + o(1)$.

Analysis of $\{\theta_n\}$ over the slow time scale: Going back to the slow time scale defined by (55), define the continuous time processes $\{\bar{w}_t, \bar{c}_t : t \geq 0\}$ as before. Define similarly the piecewise constant time processes $\{\bar{\mathcal{A}}_t, \bar{\mathcal{G}}_t : t \geq 0\}$ as well as the three components $\{\bar{\mathcal{A}}_t^+, \bar{\mathcal{A}}_t^-, \bar{\mathcal{A}}_t^\zeta : t \geq 0\}$.

Proposition A.20. *The conclusions of Prop. A.11 and Prop. A.12 hold for general eligibility vectors, subject to the modified definition of $\mathcal{A}(\theta)$ in (81).*

Proof. For the three claims of Prop. A.11, it suffices to prove that Prop. A.11 (i) holds with the new definition (81) of $\mathcal{A}(\theta)$. The rest of the claims then follow from (i).

At a point t where both \bar{w}_t and \bar{c}_t are differentiable, denote $v_t = \frac{d}{dt} \bar{w}_t$. Consider

$$\lim_{s \rightarrow 0} \frac{\bar{f}^+(\bar{w}_t, \bar{w}_{t+s}) - \bar{f}^+(\bar{w}_t, \bar{w}_t)}{s} = \lim_{s \rightarrow 0} \sum_{x, u} \varpi(x, u) \zeta^+(\bar{w}_t, x, u) \frac{\varsigma^{\bar{w}_{t+s}}(x, u) - \varsigma^{\bar{w}_t}(x, u)}{s}$$

By Lemma A.8, $\varsigma^{\bar{w}_t}(x, u)$ is differentiable for each state-action pair and a.e. t , and hence

$$\bar{f}^+(\bar{w}_t; v_t) = -\bar{f}^+(\bar{w}_t; -v_t), \quad \text{for a.e. } t \in \mathbb{R}_+$$

The same arguments imply $\bar{f}^-(\bar{w}_t; v_t) = -\bar{f}^-(\bar{w}_t; -v_t)$ for a.e. $t \in \mathbb{R}_+$. Then, with Prop. A.19 (i), the same arguments used to establish Prop. A.11 (i) yield those conclusions: For each $T > 0$,

$$\begin{aligned} \lim_{T_0 \rightarrow \infty} \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t^+ v_t - \bar{f}^+(\bar{w}_t; v_t)\|_\infty dt &= 0 \\ \lim_{T_0 \rightarrow \infty} \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t^- v_t - \bar{f}^-(\bar{w}_t; v_t)\|_\infty dt &= 0 \end{aligned}$$

It follows from (82) that

$$\frac{d}{dt} \bar{c}_t = \bar{f}^+(\bar{w}_t; v_t) - \bar{f}^-(\bar{w}_t; v_t) + \mathbb{E}_\varpi[\mathcal{D}(\bar{w}_t, \Phi_{n+1}) \partial_\xi \zeta_n] v_t$$

Therefore,

$$\begin{aligned} \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t v_t - \frac{d}{dt} \bar{c}_t\|_\infty dt &\leq \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t^+ v_t - \bar{f}^+(\bar{w}_t; v_t)\|_\infty + \|\bar{\mathcal{A}}_t^- v_t - \bar{f}^-(\bar{w}_t; v_t)\|_\infty dt \\ &\quad + \int_{T_0}^{T_0+T} \|\bar{\mathcal{A}}_t^\zeta v_t - \mathbb{E}_\varpi[\mathcal{D}(\bar{w}_t, \Phi_{n+1}) \partial_\xi \zeta_n] v_t\|_\infty dt \end{aligned}$$

where the right hand side of the above inequality goes to 0 as $T_0 \rightarrow \infty$.

For the conclusions of Prop. A.12, we only need to prove (ii) with the new $\mathcal{A}(\theta)$. Let $\mathcal{A}_t^+, \mathcal{A}_t^-, \mathcal{A}_t^\zeta$ denote the weak sub-sequential limits of $\{\bar{\mathcal{A}}_{T_0+t}^+, \bar{\mathcal{A}}_{T_0+t}^-, \bar{\mathcal{A}}_{T_0+t}^\zeta : T_0 \geq 0, 0 \leq t \leq T\}$ respectively. By Prop. A.19 (ii), the same arguments used for Prop. A.12 (ii) apply to each of \mathcal{A}_t^+ and \mathcal{A}_t^- ,

$$\mathcal{A}_t^+ \in \mathcal{A}^+(w_t), \quad \mathcal{A}_t^- \in \mathcal{A}^-(w_t), \quad t \in [0, T]$$

We also have $\mathcal{A}_t^\zeta = \mathbb{E}_\varpi[\mathcal{D}(w_t, \Phi_{n+1}) \partial_\xi \zeta_n]$ from Prop. A.19 (iii). Therefore, $\mathcal{A}_t = \mathcal{A}_t^+ - \mathcal{A}_t^- + \mathcal{A}_t^\zeta$, and $\mathcal{A}_t \in \mathcal{A}(w_t)$ for each $t \in [0, T]$. \square

Following the same arguments as in Section A.5.1, the ODE approximations and ODE limits established in Prop. A.20 imply the following extension of Thm. 2.1:

Theorem A.21. *The conclusions of Thm. 2.1 hold, subject to the modified definition of $\mathcal{A}(\theta)$ in (81).*

A.6 Numerical Results: Implementation details

Complexity of Zap Q-learning For the Zap Q-learning algorithm (21), per-iteration complexity comes from various sources:

- (i) Computation of $f(\theta_n, \Phi_{n+1})$ involves a maximum to obtain Q^{θ_n} in (19a).
- (ii) The derivatives $A_{n+1} = \partial_\theta f(\theta_n, \Phi_{n+1})$ are easily computed for linear parameterization of Q^θ , but require back-propagation in a neural network function approximation architecture.
- (iii) Computation of $G_{n+1} f(\theta_n, \Phi_{n+1})$ in (21c) and (21d) requires (i) multiplication of two $d \times d$ matrices, and (ii) multiplying a matrix inverse and a vector. Each of these two steps has worst case computational complexity $O(d^3)$.

As discussed in Section 3, the complexity in (iii) can be reduced by updating the gain only periodically, while continuously updating estimates of $A(\theta_n)$.

The complexity bound $O(Nd^3/N_d + Nd^2)$ given in Section 3 is based on gain updates performed only at integer multiples of N_d . This bound is based on the accounting (i)–(iii) above: $O(d^2)$ complexity per iteration in (21b), and $O(d^3)$ complexity for the matrix inverse (as well as the product $\hat{A}_{n+1}^\top \hat{A}_{n+1}$ appearing in (21c)).

Meta-parameters in experiments We used $\varepsilon = 10^{-6}$ in (21c) for Mountain car and Acrobot, $\varepsilon = 10^{-4}$ for Cartpole.

For the decreasing step-size rule, we used $\rho = 0.85$ and $n_0 = 100$ in (22). For constant step-size experiments, we used

$$\alpha_n \equiv \alpha, \quad \beta_n \equiv \beta = 100\alpha$$

The choice of α itself was problem specific: $\alpha = 0.002$ for the network of size 6×3 in the Mountain car example; $\alpha = 0.005$ for other experiments using constant step-size. The average reward $\mathcal{R}(\phi^{\theta_n})$ defined in (31) was estimated by running 100 independent simulations following the policy ϕ^{θ_n} . The deterministic upper bound \bar{r} was 200 for Mountain car and Acrobot, and $\bar{r} = 1000$ for Cartpole.

Q-network The input space \mathcal{U} in each of the examples is a finite set of scalars. Recall that the size of neural networks indicated in Figure 1 refers to the size of hidden layers, with the input to the network (x, u) and the output $Q^\theta(x, u)$; hence, in the Cartpole example with $(x, u) \in \mathbb{R}^5$, the network size $30 \times 24 \times 16$ corresponds to $\theta \in \mathbb{R}^d$, with $d = 1341$:

$$d = (5 + 1) * 30 + (30 + 1) * 24 + (24 + 1) * 16 + (16 + 1) = 1341$$

where each $+ 1$ accounts for a bias parameter.

Policy The theory developed in this paper assumes a randomized stationary policy for exploration. In our experiments, we apply the parameter-dependent ϵ -greedy exploration: At iteration n ,

$$U_n = \begin{cases} \phi^{\theta_n}(X_n), & \text{with probability } 1 - \epsilon \\ \text{rand}, & \text{with probability } \epsilon \end{cases}$$

We set $\epsilon = 0.4$ for the Mountain Car and Acrobot, and $\epsilon = 0.2$ for Cartpole.