

# Supplementary Material

## A Two time-scale stochastic approximation result from Borkar [4, Chapter 6]

We provide here a well-established convergence result that we use to establish our main results. The result is taken from a work on stochastic approximation [4] and is used to establish convergence of two time-scale algorithms. Let

$$\mathbf{u}_{t+1} \leftarrow \mathbf{u}_t + \alpha_t (F(\mathbf{u}_t, \mathbf{v}_t) + \mathbf{m}_{t+1}) \quad (7a)$$

$$\mathbf{v}_{t+1} \leftarrow \mathbf{v}_t + \beta_t (G(\mathbf{u}_t, \mathbf{v}_t) + \mathbf{n}_{t+1}), \quad (7b)$$

denote two coupled iterations of a stochastic approximation algorithm, where  $\mathbf{u}_t \in \mathbb{R}^p$  and  $\mathbf{v}_t \in \mathbb{R}^q$  for all  $t$ . We consider the following assumptions:

- (A)  $F : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}^p$  and  $G : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}^q$  are both Lipschitz continuous functions;
- (B) The step size sequences  $\{\alpha_t, t \in \mathbb{N}\}$  and  $\{\beta_t, t \in \mathbb{N}\}$ , are such that

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha_t &= \infty & \sum_{t=0}^{\infty} \alpha_t^2 &< \infty, \\ \sum_{t=0}^{\infty} \beta_t &= \infty & \sum_{t=0}^{\infty} \beta_t^2 &< \infty, \end{aligned}$$

and  $\alpha_t = o(\beta_t)$ .

- (C)  $\{\mathbf{m}_t, t \in \mathbb{N}\}$  and  $\{\mathbf{n}_t, t \in \mathbb{N}\}$  are two martingale difference sequences w.r.t. the  $\sigma$ -algebra  $\mathcal{F}_t$  generated by  $\{(\mathbf{u}_\tau, \mathbf{v}_\tau, \mathbf{m}_\tau, \mathbf{n}_\tau), \tau = 0, \dots, t\}$ . Furthermore, there exist constants  $c_m, c_n$  such that for all  $t \geq 0$

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{m}_{t+1}\|^2 | \mathcal{F}_t \right] &\leq c_m (1 + \|\mathbf{u}_t\|^2 + \|\mathbf{v}_t\|^2), \\ \mathbb{E} \left[ \|\mathbf{n}_{t+1}\|^2 | \mathcal{F}_t \right] &\leq c_n (1 + \|\mathbf{u}_t\|^2 + \|\mathbf{v}_t\|^2). \end{aligned}$$

We then have the following result [4].

**Theorem 3.** *Assume that, for every  $\mathbf{u} \in \mathbb{R}^p$ , the ordinary differential equation (o.d.e.)*

$$\dot{\mathbf{v}}_t = G(\mathbf{u}, \mathbf{v}_t)$$

*has a unique, globally asymptotically stable equilibrium  $\boldsymbol{\lambda}(\mathbf{u})$ , where  $\boldsymbol{\lambda} : \mathbb{R}^p \rightarrow \mathbb{R}^q$  is Lipschitz continuous. Further assume that the o.d.e.*

$$\dot{\mathbf{u}}_t = F(\mathbf{u}_t, \boldsymbol{\lambda}(\mathbf{u}_t))$$

*has a unique, globally asymptotically stable equilibrium  $\mathbf{u}^* \in \mathbb{R}^p$ . Then, under Assumptions (A) through (C), the coupled iterations (7) converge w.p.1 to  $(\mathbf{u}^*, \boldsymbol{\lambda}(\mathbf{u}^*))$  as long as  $\sup_t \|\mathbf{u}_t\| < \infty$  and  $\sup_t \|\mathbf{v}_t\| < \infty$  w.p.1.*

## B Proof of Theorem 1

This appendix provides a more detailed proof of Theorem 1, carefully establishing each of the technical conditions required for the application of Theorem 3. We establish a number of intermediate results (Propositions 1 through 4) that establish the key properties of the mean fields  $F$  and  $G$  and the martingale difference sequences  $\{\mathbf{m}_t\}$  and  $\{\mathbf{n}_t\}$  defined in the main text. We then establish, in Propositions 5 and 6, the stable behavior for the relevant o.d.e..

## B.1 Preliminaries

For convenience, we repeat herein the relevant definitions from the main text. Our algorithm is defined by the two coupled updates

$$\mathbf{u}_{t+1} \leftarrow \mathbf{u}_t + \alpha_t (\phi(x_t, a_t) Q_{\mathbf{v}_t}(x_t, a_t) - \mathbf{u}_t), \quad (4a)$$

$$\mathbf{v}_{t+1} \leftarrow \mathbf{v}_t + \beta_t \phi(x_t, a_t) \delta_t, \quad (4b)$$

where

$$\delta_t = r_t + \gamma \max_{a' \in \mathcal{A}} Q_{\mathbf{u}_t}(x'_t, a') - Q_{\mathbf{v}_t}(x_t, a_t).$$

We assume that

- (I) For all  $t$ ,  $(x_t, a_t, x'_t, r_t)$  is sampled from  $\mathcal{B} = \{(x_i, a_i, x'_i, r_i), i \in \mathbb{N}_0\}$  according to a fixed distribution  $\mu$  over  $\mathcal{B}$ . Moreover, the next-state distribution  $x'$  is such that  $\mu(x' | x, a) = \mathbf{P}(x' | x, a)$  for each  $x' \in \mathcal{X}$  and the reward distribution  $r$  is such that  $\mathbb{E}_\mu[r | x, a] = R(x, a)$ .
- (II) The matrix  $\mathbb{E}_\mu [\phi(x_t, a_t) \phi^T(x_t, a_t)]$  is non-singular and  $\|\phi(x, a)\|_2 \leq 1$ , for all pairs  $(x, a) \in \mathcal{X} \times \mathcal{A}$ .
- (III) The step size sequences  $\{\alpha_t, t \in \mathbb{N}\}$  and  $\{\beta_t, t \in \mathbb{N}\}$ , verify

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha_t &= \infty, & \sum_{t=0}^{\infty} \alpha_t^2 &< \infty, \\ \sum_{t=0}^{\infty} \beta_t &= \infty, & \sum_{t=0}^{\infty} \beta_t^2 &< \infty, \end{aligned}$$

and, moreover,  $\alpha_t = o(\beta_t)$ .

We define the mean fields

$$F(\mathbf{u}_t, \mathbf{v}_t) \stackrel{\text{def}}{=} \mathbb{E}_\mu [\phi(x_t, a_t) \phi^T(x_t, a_t)] \mathbf{v}_t - \mathbf{u}_t, \quad (8a)$$

$$G(\mathbf{u}_t, \mathbf{v}_t) \stackrel{\text{def}}{=} \mathbb{E}_\mu [\phi(x_t, a_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t]. \quad (8b)$$

and the martingale differences

$$\mathbf{m}_{t+1} \stackrel{\text{def}}{=} \left( \phi(x_t, a_t) \phi^T(x_t, a_t) \mathbf{v}_t - \mathbf{u}_t \right) - F(\mathbf{u}_t, \mathbf{v}_t), \quad (9a)$$

$$\mathbf{n}_{t+1} \stackrel{\text{def}}{=} \phi(x_t, a_t) \delta_t - G(\mathbf{u}_t, \mathbf{v}_t). \quad (9b)$$

Finally, we consider the  $\sigma$ -algebra

$$\mathcal{F}_t \stackrel{\text{def}}{=} \sigma(\{\mathbf{u}_\tau, \mathbf{v}_\tau, \mathbf{m}_\tau, \mathbf{n}_\tau, \tau = 0, \dots, t\}).$$

For the upcoming results, it is important to emphasize that, from Assumption (I), the sequence  $(x_t, a_t, r_t, x'_t)$  is i.i.d., generated by a fixed distribution  $\mu$  and thus independent from  $\mathcal{F}_t$  itself. Unless otherwise noted,  $\|\cdot\|$  refers to the standard 2-norm in  $\mathbb{R}^K$ .

## B.2 Lipschitz continuity of $F$ and $G$

We start by establishing  $F$  to be Lipschitz continuous.

**Proposition 1.** *The function  $F : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}^K$ , defined in (8a), is Lipschitz continuous.*

*Proof.* We constructively show that for some  $c_F \geq 0$ , and for all  $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{z} \in \mathbb{R}^K$ ,

$$\|F(\mathbf{u}, \mathbf{v}) - F(\mathbf{w}, \mathbf{z})\| \leq c_F \|(\mathbf{u}, \mathbf{v}) - (\mathbf{w}, \mathbf{z})\|.$$

We have that

$$\begin{aligned} \|F(\mathbf{u}, \mathbf{v}) - F(\mathbf{w}, \mathbf{z})\| &\leq \left\| \mathbb{E}_\mu [\phi(x_t, a_t) \phi^T(x_t, a_t)] (\mathbf{v} - \mathbf{z}) \right\| + \|\mathbf{u} - \mathbf{w}\| \\ &\leq \left\| \mathbb{E}_\mu [\phi(x_t, a_t) \phi^T(x_t, a_t)] \right\| \|\mathbf{v} - \mathbf{z}\| + \|\mathbf{u} - \mathbf{w}\| \end{aligned}$$

where the first inequality follows from the triangle inequality, and the second follows from the Cauchy-Schwarz inequality. Using Jensen's inequality and the fact that  $\|\phi(x, a)\|_2 \leq 1$ , we have that

$$\begin{aligned} \left\| \mathbb{E}_\mu \left[ \phi(x_t, a_t) \phi^T(x_t, a_t) \right] \right\| &\leq \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \phi^T(x_t, a_t) \right\| \right] \\ &= \mathbb{E}_\mu \left[ \sup_{\|\mathbf{x}\|=1} \left\| \phi(x_t, a_t) \phi^T(x_t, a_t) \mathbf{x} \right\| \right] \\ &\leq \mathbb{E}_\mu \left[ \sup_{\|\mathbf{x}\|=1} \left\| \phi(x_t, a_t) \right\| \left\| \phi^T(x_t, a_t) \mathbf{x} \right\| \right] \\ &\leq 1. \end{aligned}$$

This yields

$$\|F(\mathbf{u}, \mathbf{v}) - F(\mathbf{w}, \mathbf{z})\| \leq \|\mathbf{v} - \mathbf{z}\| + \|\mathbf{u} - \mathbf{w}\| \leq \sqrt{K} \|(\mathbf{u}, \mathbf{v}) - (\mathbf{w}, \mathbf{z})\|,$$

and the proof is complete.  $\square$

We now establish a similar result for  $G$ .

**Proposition 2.** *The function  $G : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}^K$ , defined in (8b), is Lipschitz continuous.*

*Proof.* We want to show that, for some  $c_G \geq 0$ ,

$$\|G(\mathbf{u}, \mathbf{v}) - G(\mathbf{w}, \mathbf{z})\|_2 \leq c_G \|(\mathbf{u}, \mathbf{v}) - (\mathbf{w}, \mathbf{z})\|_2,$$

for any fixed  $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{z} \in \mathbb{R}^K$ . Following along the lines of the proof of Proposition 1, we get

$$\begin{aligned} &\|G(\mathbf{u}, \mathbf{v}) - G(\mathbf{w}, \mathbf{z})\| \\ &\leq \left\| \mathbb{E}_\mu \left[ \gamma \phi(x_t, a_t) (\max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{w}) \right] \right\| + \left\| \mathbb{E}_\mu \left[ \phi(x_t, a_t) \phi^T(x_t, a_t) (\mathbf{v} - \mathbf{z}) \right] \right\| \\ &\leq \gamma \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \right\| \left| \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{w} \right| \right] + \left\| \mathbb{E}_\mu \left[ \phi(x_t, a_t) \phi^T(x_t, a_t) \right] \right\| \|\mathbf{v} - \mathbf{z}\| \\ &\leq \gamma \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \right\| \max_{a' \in \mathcal{A}} \left| \phi^T(x'_t, a') (\mathbf{u} - \mathbf{w}) \right| \right] + \|\mathbf{v} - \mathbf{z}\| \\ &\leq \gamma \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \right\| \max_{a' \in \mathcal{A}} \left\| \phi^T(x'_t, a') \right\| \|\mathbf{u} - \mathbf{w}\| \right] + \|\mathbf{v} - \mathbf{z}\| \\ &\leq \gamma \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \right\| \max_{a' \in \mathcal{A}} \left\| \phi^T(x'_t, a') \right\| \right] \|\mathbf{u} - \mathbf{w}\| + \|\mathbf{v} - \mathbf{z}\| \\ &\leq \gamma \|(\mathbf{u}, \mathbf{v}) - (\mathbf{w}, \mathbf{z})\|, \end{aligned}$$

and the proof is complete.  $\square$

### B.3 Square integrability of $\{\mathbf{m}_t\}$ and $\{\mathbf{n}_t\}$

We start with the following preliminary result.

**Lemma 4.** *The sequence  $\{(\mathbf{m}_t, \mathcal{F}_t), t \in \mathbb{N}\}$  defined in (9a) is a martingale difference sequence.*

*Proof.* To verify that  $\{(\mathbf{m}_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is a martingale difference sequence, we must verify that the following conditions hold for every  $t$ :

- $\mathbb{E} [\|\mathbf{m}_t\|] < \infty$ ;
- $\mathbb{E} [\mathbf{m}_{t+1} \mid \mathcal{F}_t] = 0$ .

For the first bullet, we have that

$$\mathbb{E} [\|\mathbf{m}_{t+1}\|] = \mathbb{E} \left[ \left\| \left( \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) - \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \right) \mathbf{v}_t \right\| \right].^1$$

Using the Cauchy-Schwartz inequality and the linearity of the expectation, we get

$$\begin{aligned} \mathbb{E} [\|\mathbf{m}_{t+1}\|] &\leq \mathbb{E} \left[ \left\| \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) - \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \right\| \right] \mathbb{E} [\|\mathbf{v}_t\|] \\ &\leq \mathbb{E} \left[ \left\| \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right\| \right] + \left\| \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \right\| \mathbb{E} [\|\mathbf{v}_t\|]. \end{aligned}$$

Repeating the steps from the proof of Proposition 1, it follows that

$$\mathbb{E} \left[ \left\| \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right\| \right] \leq 1,$$

yielding

$$\mathbb{E} [\|\mathbf{m}_{t+1}\|] \leq 2\mathbb{E} [\|\mathbf{v}_t\|] < \infty,$$

since  $\mathbf{v}_t$  is constructed by a finite number of algebraic operations over finite quantities.<sup>2</sup>

To finish the proof, it remains to show that  $\mathbb{E} [\mathbf{m}_{t+1} | \mathcal{F}_t] = 0$ . We have that

$$\begin{aligned} \mathbb{E} [\mathbf{m}_{t+1} | \mathcal{F}_t] &= \mathbb{E} \left[ \left( \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) - \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \right) \mathbf{v}_t \mid \mathcal{F}_t \right] \\ &= \mathbb{E} \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) - \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \mid \mathcal{F}_t \right] \mathbf{v}_t \\ &= \left( \mathbb{E} \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \mid \mathcal{F}_t \right] - \mathbb{E} \left[ \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \mid \mathcal{F}_t \right] \right) \mathbf{v}_t \\ &= \left( \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] - \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) \phi^T(x_t, \mathbf{a}_t) \right] \right) \mathbf{v}_t \\ &= 0. \end{aligned}$$

□

For the sequence  $\{\mathbf{n}_t\}$ , we have a similar result.

**Lemma 5.** *The sequence  $\{(\mathbf{n}_t, \mathcal{F}_t), t \in \mathbb{N}\}$  defined in (9b) is a martingale difference sequence.*

*Proof.* As in the proof of Lemma 4, the following must hold:

- $\mathbb{E} [\|\mathbf{n}_t\|] < \infty$ ;
- $\mathbb{E} [\mathbf{n}_{t+1} | \mathcal{F}_t] = 0$ .

Writing down the expression for  $\mathbf{n}_t$ , we get

$$\begin{aligned} \mathbb{E} [\|\mathbf{n}_{t+1}\|] &= \mathbb{E} [\|\phi(x_t, \mathbf{a}_t) \delta_t - \mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t]\|] \\ &\leq \mathbb{E} [\|\phi(x_t, \mathbf{a}_t) \delta_t\| + \|\mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t]\|], \end{aligned}$$

where the inequality follows from the triangle inequality. Using Jensen's inequality, we then get

$$\mathbb{E} [\|\mathbf{n}_{t+1}\|] \leq \mathbb{E} [\|\phi(x_t, \mathbf{a}_t) \delta_t\| + \mathbb{E}_\mu [\|\phi(x_t, \mathbf{a}_t) \delta_t\| | \mathbf{u}_t, \mathbf{v}_t]] < \infty,$$

where the last inequality follows from the fact that all terms are bounded.

To show that  $\mathbb{E} [\mathbf{n}_{t+1} | \mathcal{F}_t] = 0$ , we can write

$$\begin{aligned} \mathbb{E} [\mathbf{n}_{t+1} | \mathcal{F}_t] &= \mathbb{E} [\phi(x_t, \mathbf{a}_t) \delta_t - \mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t] | \mathcal{F}_t] \\ &= \mathbb{E} [\phi(x_t, \mathbf{a}_t) \delta_t | \mathcal{F}_t] - \mathbb{E} [\mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t] | \mathcal{F}_t] \\ &= \mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t] - \mathbb{E}_\mu [\phi(x_t, \mathbf{a}_t) \delta_t | \mathbf{u}_t, \mathbf{v}_t] \\ &= 0. \end{aligned}$$

□

<sup>1</sup>Since we are not conditioning on  $\mathcal{F}_t$ ,  $\mathbf{v}_t$  is treated as a random variable.

<sup>2</sup>Note that this does not imply, however, that  $\sup_t \|\mathbf{v}_t\| < \infty$ .

Endowed with Lemmas 4 and 5, we proceed to establishing square-integrability of  $\{\mathbf{m}_t\}$  and  $\{\mathbf{n}_t\}$ .

**Proposition 3.** *There exists  $c_{\mathbf{m}} > 0$  such that, for any  $t > 0$ ,*

$$\mathbb{E} \left[ \|\mathbf{m}_{t+1}\|^2 \mid \mathcal{F}_t \right] \leq c_{\mathbf{m}}(1 + \|\mathbf{u}_t\|^2 + \|\mathbf{v}_t\|^2).$$

*Proof.* From the definition, we get

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{m}_{t+1}\|^2 \mid \mathcal{F}_t \right] &= \mathbb{E} \left[ \left\| \left( \phi(x_t, a_t) \phi^T(x_t, a_t) - \mathbb{E}_\mu \left[ \phi(x_t, a_t) \phi^T(x_t, a_t) \right] \right) \mathbf{v}_t \right\|^2 \mid \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \left( \left\| \phi(x_t, a_t) \phi^T(x_t, a_t) \right\| + \left\| \mathbb{E}_\mu \left[ \phi(x_t, a_t) \phi^T(x_t, a_t) \right] \right\| \right)^2 \|\mathbf{v}_t\|^2 \mid \mathcal{F}_t \right] \\ &\leq 4\|\mathbf{v}_t\|^2 \\ &\leq 4(1 + \|\mathbf{u}_t\|_2^2 + \|\mathbf{v}_t\|_2^2). \end{aligned}$$

□

For  $\{\mathbf{n}_t\}$ , we get a similar result.

**Proposition 4.** *There exists  $c_{\mathbf{n}} > 0$  such that, for any  $t > 0$ ,*

$$\mathbb{E} \left[ \|\mathbf{n}_{t+1}\|^2 \mid \mathcal{F}_t \right] \leq c_{\mathbf{n}}(1 + \|\mathbf{u}_t\|^2 + \|\mathbf{v}_t\|^2).$$

*Proof.* Note that

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{n}_{t+1}\|^2 \mid \mathcal{F}_t \right] &= \mathbb{E} \left[ \left\| \phi(x_t, a_t) \delta_t - \mathbb{E}_\mu \left[ \phi(x_t, a_t) \delta_t \mid \mathbf{u}_t, \mathbf{v}_t \right] \right\|^2 \mid \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \left( \|\phi(x_t, a_t) \delta_t\| + \left\| \mathbb{E} \left[ \phi(x_t, a_t) \delta_t \mid \mathbf{u}_t, \mathbf{v}_t \right] \right\| \right)^2 \mid \mathcal{F}_t \right] \\ &= 4\mathbb{E} \left[ \|\phi(x_t, a_t) \delta_t\|^2 \mid \mathcal{F}_t \right]. \end{aligned}$$

Breaking down the right-hand side, we get

$$\begin{aligned} \mathbb{E} \left[ \|\phi(x_t, a_t) \delta_t\|^2 \mid \mathcal{F}_t \right] &= \mathbb{E} \left[ \left\| \phi(x_t, a_t) \left( r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}_t - \phi^T(x_t, a_t) \mathbf{v}_t \right) \right\|^2 \mid \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \|\phi(x_t, a_t)\|^2 \left| r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}_t - \phi^T(x_t, a_t) \mathbf{v}_t \right|^2 \mid \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \rho + \left( \left| \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}_t \right| + \left| \phi^T(x_t, a_t) \mathbf{v}_t \right| \right)^2 \mid \mathcal{F}_t \right] \\ &\leq \mathbb{E} \left[ \rho + (\gamma \max_{a' \in \mathcal{A}} \|\phi^T(x_t, a')\| \|\mathbf{u}_t\| + \|\phi^T(x_t, a_t)\| \|\mathbf{v}_t\|)^2 \mid \mathcal{F}_t \right] \\ &\leq \rho + (\|\mathbf{u}_t\| + \|\mathbf{v}_t\|)^2 \\ &\leq c_{\mathbf{n}}(1 + \|\mathbf{u}_t\|^2 + \|\mathbf{v}_t\|^2), \end{aligned}$$

where  $c_{\mathbf{n}}$  depends on  $\rho$  and  $K$ .

□

#### B.4 Stability of the o.d.e.

Propositions 1 through 4 establish that our algorithm verifies the technical assumptions of Theorem 3. It remains to show that the remaining conditions of theorem also hold—namely, the stability of the associated o.d.e. and the boundedness of the iterates  $\mathbf{u}_t$  and  $\mathbf{v}_t$ .

**Proposition 5.** *For any fixed  $\mathbf{u} \in \mathbb{R}^K$ , the ordinary differential equation*

$$\dot{\mathbf{v}}_t = G(\mathbf{u}, \mathbf{v}_t) \tag{10}$$

*has a unique, globally asymptotically stable equilibrium  $\boldsymbol{\lambda}(\mathbf{u})$ , where  $\boldsymbol{\lambda} : \mathbb{R}^K \rightarrow \mathbb{R}^K$  is Lipschitz continuous.*

*Proof.* For a given  $\mathbf{u}$ ,  $\mathbf{v}^* \in \mathbb{R}^K$  is an equilibrium of the o.d.e. (10) if

$$G(\mathbf{u}, \mathbf{v}^*) = 0$$

or, equivalently, if

$$\mathbf{v}^* = \Sigma_\mu^{-1} \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}) \right].$$

That such equilibrium exists for any  $\mathbf{u}$  follows from the fact that  $\Sigma_\mu$  is non-singular, as required by Assumption (II). Define  $\boldsymbol{\lambda} : \mathbb{R}^K \rightarrow \mathbb{R}^K$  as

$$\boldsymbol{\lambda}(\mathbf{u}) = \Sigma_\mu^{-1} \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}) \right].$$

To see that the function  $\boldsymbol{\lambda}$  thus defined is Lipschitz continuous, we note that

$$\begin{aligned} \|\boldsymbol{\lambda}(\mathbf{u}) - \boldsymbol{\lambda}(\mathbf{u}')\| &= \left\| \Sigma_\mu^{-1} \mathbb{E}_\mu \left[ \gamma \phi(x_t, \mathbf{a}_t) (\max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}') \right] \right\| \\ &\leq \left\| \Sigma_\mu^{-1} \right\| \left\| \mathbb{E}_\mu \left[ \gamma \phi(x_t, \mathbf{a}_t) (\max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}') \right] \right\|. \end{aligned}$$

Using Jensen's inequality, we get

$$\begin{aligned} \|\boldsymbol{\lambda}(\mathbf{u}) - \boldsymbol{\lambda}(\mathbf{u}')\| &\leq \left\| \Sigma_\mu^{-1} \right\| \mathbb{E}_\mu \left[ \left\| \gamma \phi(x_t, \mathbf{a}_t) (\max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}') \right\| \right] \\ &\leq \left\| \Sigma_\mu^{-1} \right\| \mathbb{E}_\mu \left[ \gamma \|\phi(x_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \|\phi^T(x'_t, a')(\mathbf{u} - \mathbf{u}')\| \right] \\ &\leq \left\| \Sigma_\mu^{-1} \right\| \mathbb{E}_\mu \left[ \gamma \|\phi(x_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \|\phi^T(x'_t, a')\| \|\mathbf{u} - \mathbf{u}'\| \right] \\ &\leq c_\lambda \|\mathbf{u} - \mathbf{u}'\|, \end{aligned}$$

for some  $c_\lambda > 0$ .

Finally, to show that, given  $\mathbf{u}$ ,  $\boldsymbol{\lambda}(\mathbf{u})$  is a globally asymptotically stable equilibrium for the o.d.e. (10), we define the candidate Lyapunov function  $L_G : \mathbb{R}^K \rightarrow \mathbb{R}$  as

$$L_G(\mathbf{v}) = \frac{1}{2} \|\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u})\|^2.$$

The proof is complete as long as

1.  $L_G(\mathbf{v}) \geq 0$  for all  $\mathbf{v} \in \mathbb{R}^K$ ;
2.  $L_G(\mathbf{v}) = 0$  if and only if  $\mathbf{v} = \boldsymbol{\lambda}(\mathbf{u})$ ;
3.  $\dot{L}_G(\mathbf{v}) \leq 0$  for all  $\mathbf{v} \in \mathbb{R}^K$ ;
4.  $\dot{L}_G(\mathbf{v}) = 0$  if and only if  $\mathbf{v} = \boldsymbol{\lambda}(\mathbf{u})$ .

The first two properties follow directly from the definition of  $L_G$ . As for the last two, we start by explicitly writing  $\dot{L}_G$ , to get

$$\dot{L}_G(\mathbf{v}) \stackrel{\text{def}}{=} \frac{d}{dt} L_G(\mathbf{v}) = \sum_{k=1}^K \frac{\partial L_G}{\partial v_k} G_k(\mathbf{u}, \mathbf{v}) = (\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u}))^T G(\mathbf{u}, \mathbf{v})$$

Hence,

$$\begin{aligned} \dot{L}_G(\mathbf{v}) &= (\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u}))^T \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \phi^T(x_t, \mathbf{a}_t) \mathbf{v}) \right] \\ &= (\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u}))^T \left( \mathbb{E}_\mu \left[ \phi(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \phi^T(x_t, \mathbf{a}_t) \mathbf{v}) \right] - G(\mathbf{u}, \boldsymbol{\lambda}(\mathbf{u})) \right), \end{aligned}$$

where we used the fact that  $G(\mathbf{u}, \boldsymbol{\lambda}(\mathbf{u})) = 0$ . Then, after some shuffling, we finally get

$$\begin{aligned} &= (\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u}))^T \mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) \boldsymbol{\phi}^T(x_t, \mathbf{a}_t) (\boldsymbol{\lambda}(\mathbf{u}) - \mathbf{v}) \right] \\ &= -(\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u}))^T \boldsymbol{\Sigma}_\mu (\mathbf{v} - \boldsymbol{\lambda}(\mathbf{u})) \leq 0, \end{aligned}$$

where the last inequality comes from the fact that  $\mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) \boldsymbol{\phi}^T(x_t, \mathbf{a}_t) \right]$  is an auto-covariance matrix and, as such, positive definite. The conclusion follows.  $\square$

Next, we present a similar result for the slower o.d.e.

**Proposition 6.** *The ordinary differential equation*

$$\dot{\mathbf{u}}_t = F(\mathbf{u}_t, \boldsymbol{\lambda}(\mathbf{u}_t)) \tag{11}$$

has a unique, globally asymptotically stable equilibrium  $\mathbf{u}^* \in \mathbb{R}^K$ .

*Proof.* We start by establishing the existence of at least one equilibrium. A vector  $\mathbf{u}^* \in \mathbb{R}^K$  is an equilibrium for the o.d.e. (11) if

$$F(\mathbf{u}^*, \boldsymbol{\lambda}(\mathbf{u}^*)) = 0$$

or, equivalently, if

$$\mathbf{u}^* = \mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{u}^*) \right].$$

Define  $H' : \mathbb{R}^K \rightarrow \mathbb{R}^K$  as

$$H'(\mathbf{u}) = \mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{u}) \right].$$

Then, the equilibria for the o.d.e. (11) correspond to the fixed points of  $H'$ . By showing  $H'$  to be a contraction, Banach's fixed point theorem ensures existence of a single fixed point, thus establishing both the existence and unicity of an equilibrium for (11). Given any  $\mathbf{z}, \mathbf{w} \in \mathbb{R}^K$ ,

$$\begin{aligned} &\|H'(\mathbf{w}) - H'(\mathbf{z})\| \\ &= \left\| \mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{w}) \right] - \mathbb{E}_\mu \left[ \boldsymbol{\phi}(x_t, \mathbf{a}_t) (r_t + \gamma \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{z}) \right] \right\| \\ &= \left\| \mathbb{E}_\mu \left[ \gamma \boldsymbol{\phi}(x_t, \mathbf{a}_t) (\max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{w} - \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{z}) \right] \right\|. \end{aligned}$$

Using Jensen's inequality,

$$\begin{aligned} \|H(\mathbf{w}) - H(\mathbf{z})\| &\leq \mathbb{E}_\mu \left[ \left\| \gamma \boldsymbol{\phi}(x_t, \mathbf{a}_t) (\max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{w} - \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{z}) \right\| \right] \\ &= \mathbb{E}_\mu \left[ \gamma \|\boldsymbol{\phi}(x_t, \mathbf{a}_t)\| \left| \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{w} - \max_{a' \in \mathcal{A}} \boldsymbol{\phi}^T(x'_t, a') \mathbf{z} \right| \right] \\ &\leq \mathbb{E}_\mu \left[ \gamma \|\boldsymbol{\phi}(x_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left| \boldsymbol{\phi}^T(x'_t, a') \mathbf{w} - \boldsymbol{\phi}^T(x'_t, a') \mathbf{z} \right| \right] \\ &= \mathbb{E}_\mu \left[ \gamma \|\boldsymbol{\phi}(x_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left| \boldsymbol{\phi}^T(x'_t, a') (\mathbf{w} - \mathbf{z}) \right| \right] \\ &\leq \mathbb{E}_\mu \left[ \gamma \|\boldsymbol{\phi}(x_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left\| \boldsymbol{\phi}^T(x'_t, a') \right\| \|\mathbf{w} - \mathbf{z}\| \right] \\ &\leq \gamma \|\mathbf{w} - \mathbf{z}\|. \end{aligned}$$

It follows that there is, in fact, a unique equilibrium  $\mathbf{u}^*$  for the o.d.e. To show that it is globally asymptotically stable, we again use Lyapunov's second method. We define the candidate Lyapunov function  $L_F : \mathbb{R}^K \rightarrow \mathbb{R}$  as

$$L_F(\mathbf{u}) = \frac{1}{2} \|\mathbf{u} - \mathbf{u}^*\|^2.$$

Once again, the conclusion follows as long as

1.  $L_F(\mathbf{u}) \geq 0$  for all  $\mathbf{u} \in \mathbb{R}^K$ ;
2.  $L_F(\mathbf{u}) = 0$  if and only if  $\mathbf{u} = \mathbf{u}^*$ ;
3.  $\dot{L}_F(\mathbf{u}) \leq 0$  for all  $\mathbf{u} \in \mathbb{R}^K$ ;
4.  $\dot{L}_F(\mathbf{u}) = 0$  if and only if  $\mathbf{u} = \mathbf{u}^*$ .

The first two follow directly from the definition of  $L_F$ . As for the last two,

$$\begin{aligned} \dot{L}_F(\mathbf{u}) &= (\mathbf{u} - \mathbf{u}^*)^T F(\mathbf{u}, \boldsymbol{\lambda}(\mathbf{u}))^T \\ &= (\mathbf{u} - \mathbf{u}^*)^T \mathbb{E}_\mu \left[ \gamma \phi(x_t, a_t) \left( \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}^* \right) \right] - \|\mathbf{u} - \mathbf{u}^*\|^2, \end{aligned}$$

where we used the fact that  $F(\mathbf{u}^*, \boldsymbol{\lambda}(\mathbf{u}^*)) = 0$ . Using Jensen's inequality, we get

$$\begin{aligned} \dot{L}_F(\mathbf{u}) &\leq \|\mathbf{u} - \mathbf{u}^*\| \mathbb{E}_\mu \left[ \left\| \gamma \phi(x_t, a_t) \left( \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u} - \max_{a' \in \mathcal{A}} \phi^T(x'_t, a') \mathbf{u}^* \right) \right\| \right] - \|\mathbf{u} - \mathbf{u}^*\|^2 \\ &\leq \gamma \mathbb{E}_\mu \left[ \left\| \phi(x_t, a_t) \right\| \max_{a' \in \mathcal{A}} \left\| \phi^T(x'_t, a') \right\| \right] \|\mathbf{u} - \mathbf{u}^*\|^2 - \|\mathbf{u} - \mathbf{u}^*\|^2 \\ &\leq \gamma \|\mathbf{u} - \mathbf{u}^*\|^2 - \|\mathbf{u} - \mathbf{u}^*\|^2 \\ &\leq (\gamma - 1) \|\mathbf{u} - \mathbf{u}^*\|^2 \leq 0. \end{aligned}$$

The conclusion follows. □

## C Boundedness of $\mathbf{u}_t$ and $\mathbf{v}_t$

We conclude by establishing the boundedness of the iterates  $\mathbf{v}_t$  and  $\mathbf{u}_t$  under Assumptions (I) through (III). To do that, we use the following result.

**Theorem 6** ([5]). *Given the algorithm*

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t (H(\mathbf{w}_t) + \mathbf{m}_{t+1}),$$

where

1. The function  $H : \mathbb{R}^K \rightarrow \mathbb{R}^K$  is Lipschitz continuous. Moreover, defining  $H_r : \mathbb{R}^K \rightarrow \mathbb{R}^K$  as

$$H_r(\mathbf{w}) = \frac{H(r\mathbf{w})}{r},$$

there is a function  $H_\infty : \mathbb{R}^K \rightarrow \mathbb{R}^K$  such that

$$\lim_{r \rightarrow \infty} H_r(\mathbf{w}) = H_\infty(\mathbf{w})$$

for all  $\mathbf{w} \in \mathbb{R}^K$ .

2. The origin is a globally asymptotically stable equilibrium of the o.d.e.

$$\dot{\mathbf{w}}_t = H_\infty(\mathbf{w}_t).$$

3. The sequence  $\{\mathbf{m}_t, t \in \mathbb{N}\}$  is a martingale difference sequence and verifies, for all  $t \geq 0$

$$\mathbb{E} \left[ \|\mathbf{m}_{t+1}\|^2 \mid \mathcal{F}_t \right] \leq c_0 (1 + \|\mathbf{w}_t\|^2)$$

for some  $c_0 > 0$ .

4. The sequence  $\{\alpha_t, t \in \mathbb{N}\}$  verifies

$$\sum_{t=0}^{\infty} \alpha_t = \infty, \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty.$$

Then, with probability 1,  $\sup_t \|\mathbf{w}_t\| < \infty$ .

Following [4], we analyze each iterate of the algorithm separately. In particular, we analyze the faster iteration by treating the slower as stationary and analyze the slower iteration by treating the faster as if in equilibrium.

For a fixed  $\mathbf{u} \in \mathbb{R}^K$ , let

$$G_c(\mathbf{v}) = \frac{G(\mathbf{u}, c\mathbf{v})}{c}.$$

We have the following result.

**Lemma 7.** *There is a limiting function  $G_\infty : \mathbb{R}^K \rightarrow \mathbb{R}^K$  such that*

$$\lim_{c \rightarrow \infty} G_c(\mathbf{v}) = G_\infty(\mathbf{v}),$$

*and the origin is an asymptotically stable equilibrium for the o.d.e*

$$\dot{\mathbf{v}}_t = G_\infty(\mathbf{v}_t).$$

*Proof.* Replacing the definition of  $G_c$ , we get

$$\begin{aligned} \lim_{c \rightarrow \infty} G_c(\mathbf{v}) &= \lim_{c \rightarrow \infty} \frac{G(\mathbf{u}, c\mathbf{v})}{c} \\ &= \lim_{c \rightarrow \infty} \frac{1}{c} \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \left( r_t + \gamma \max_{\mathbf{a}' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, \mathbf{a}_t) \mathbf{u} - c \phi^T(\mathbf{x}_t, \mathbf{a}_t) \mathbf{v} \right) \right] \\ &= -\Sigma_\mu \mathbf{v}. \end{aligned}$$

Therefore, letting  $G_\infty(\mathbf{v}) = -\Sigma_\mu \mathbf{v}$ , the o.d.e.

$$\dot{\mathbf{v}}_t = G_\infty(\mathbf{v}_t) = -\Sigma_\mu \mathbf{v}_t.$$

is linear and time-invariant. Since  $\Sigma_\mu$  is positive definite, it is immediate that the origin is a globally asymptotically stable equilibrium.  $\square$

Noting that all other conditions follow directly from Assumptions (I) through (III) and Propositions 1 through 6, we can now apply Theorem 6 to get the following conclusion.

**Proposition 7.** *Under assumptions (I) through (III),  $\sup_t \|\mathbf{v}_t\| < \infty$  almost surely.*

For the slower iterate, we consider that the faster has converged, and define

$$F_c(\mathbf{u}) = \frac{F(c\mathbf{u}, \lambda(c\mathbf{u}))}{c}.$$

We have the following counterpart to Lemma 7.

**Lemma 8.** *There is a limiting function  $F_\infty : \mathbb{R}^K \rightarrow \mathbb{R}^K$  such that*

$$\lim_{c \rightarrow \infty} F_c(\mathbf{u}) = F_\infty(\mathbf{u}),$$

*and the origin is an asymptotically stable equilibrium for the o.d.e*

$$\dot{\mathbf{u}}_t = F_\infty(\mathbf{u}_t).$$

*Proof.* Replacing the definition of  $F_c$ , we get

$$\begin{aligned} \lim_{c \rightarrow \infty} F_c(\mathbf{u}) &= \lim_{c \rightarrow \infty} \frac{1}{c} \left( \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \phi^T(\mathbf{x}_t, \mathbf{a}_t) \right] \lambda(c\mathbf{u}) - c\mathbf{u} \right) \\ &= \gamma \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \max_{\mathbf{a}' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, \mathbf{a}') \mathbf{u} \right] - \mathbf{u}. \end{aligned}$$

Let us then define

$$F_\infty(\mathbf{u}) = \gamma \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \max_{\mathbf{a}' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, \mathbf{a}') \mathbf{u} \right] - \mathbf{u}.$$

Consider the candidate Lyapunov function  $L(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}\|_2^2$ . The conclusion of the lemma follows by showing that

1.  $L(\mathbf{u}) \geq 0$  for all  $\mathbf{u} \in \mathbb{R}^K$ ;
2.  $L(\mathbf{u}) = 0$  if and only if  $\mathbf{u} = \mathbf{0}$ ;
3.  $\dot{L}(\mathbf{u}) \leq 0$  for all  $\mathbf{u} \in \mathbb{R}^K$ ;
4.  $\dot{L}(\mathbf{u}) = 0$  if and only if  $\mathbf{u} = \mathbf{0}$ .

The first two conditions follow directly from the definition of  $L$ . As for the last two, we observe that

$$\begin{aligned}
\dot{L}(\mathbf{u}) &= \mathbf{u}^T F_\infty(\mathbf{u}) \\
&= \mathbf{u}^T \left( \gamma \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \max_{a' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, a') \mathbf{u} \right] - \mathbf{u} \right) \\
&= \gamma \mathbf{u}^T \mathbb{E}_\mu \left[ \phi(\mathbf{x}_t, \mathbf{a}_t) \max_{a' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, a') \mathbf{u} \right] - \|\mathbf{u}\|^2.
\end{aligned}$$

Using Jensen's inequality, we get

$$\begin{aligned}
\dot{L}(\mathbf{u}) &\leq \gamma \|\mathbf{u}\| \mathbb{E}_\mu \left[ \left\| \phi(\mathbf{x}_t, \mathbf{a}_t) \max_{a' \in \mathcal{A}} \phi^T(\mathbf{x}'_t, a') \mathbf{u} \right\| \right] - \|\mathbf{u}\|^2 \\
&\leq \gamma \|\mathbf{u}\| \mathbb{E}_\mu \left[ \|\phi(\mathbf{x}_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left| \phi^T(\mathbf{x}'_t, a') \mathbf{u} \right| \right] - \|\mathbf{u}\|^2 \\
&\leq \gamma \|\mathbf{u}\| \mathbb{E}_\mu \left[ \|\phi(\mathbf{x}_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left| \phi^T(\mathbf{x}'_t, a') \right| \right] - \|\mathbf{u}\|^2 \\
&\leq \gamma \|\mathbf{u}\| \mathbb{E}_\mu \left[ \|\phi(\mathbf{x}_t, \mathbf{a}_t)\| \max_{a' \in \mathcal{A}} \left\| \phi^T(\mathbf{x}'_t, a') \right\| \|\mathbf{u}\| \right] - \|\mathbf{u}\|^2 \\
&\leq (\gamma - 1) \|\mathbf{u}\|^2 \leq 0.
\end{aligned}$$

The proof is complete. □

Finally, again resorting to Theorem 6, we can state the following proposition.

**Proposition 8.**  $\sup_t \|\mathbf{u}_t\|_1 < \infty$  almost surely.