

Linear Dynamics meets Linear MDPs: Closed-Form Optimal Policies via Reinforcement Learning

Anonymous authors
Paper under double-blind review

Abstract

Many applications—including power systems, robotics, and economics—involve a dynamical system interacting with a stochastic and hard-to-model environment. We adopt a reinforcement learning approach to control such systems. Specifically, we consider a deterministic, discrete-time, linear, time-invariant dynamical system coupled with a feature-based linear Markov process with an unknown transition kernel. The objective is to learn a control policy that minimizes a quadratic cost over the system state, the Markov process, and the control input. Leveraging both components of the system, we derive an explicit parametric form for the optimal state-action value function and the corresponding optimal policy. Our model is distinct in combining aspects of both classical Linear Quadratic Regulator (LQR) and linear Markov decision process (MDP) frameworks. This combination retains the implementation simplicity of LQR, while allowing for sophisticated stochastic modeling afforded by linear MDPs, without estimating the transition probabilities, thereby enabling direct policy improvement. For the nominal setting, where the linear system dynamics are known, we use tools from control theory to provide theoretical guarantees on the stability of the system under the learned policy and provide a sample complexity analysis for its convergence to the optimal policy. We further extend our framework to systems with Gaussian process noise and to systems with unknown linear dynamics. We illustrate our results via numerical examples for the nominal, noisy, and unknown dynamics settings to demonstrate the effectiveness of our approach in learning the optimal control policy under partially known stochastic dynamics.

1 Introduction

In many applications, a well-modeled agent must interact with and make decisions in stochastic and hard-to-model environments, with the aim to optimize a certain cost that is affected by the agent’s objective and the environment. A prominent example arises in power systems, where a controllable energy storage device evolves under known physical dynamics, yet must respond to uncertain net load demand driven by exogenous factors such as variability in generation, consumer behavior, and weather conditions, all of which are unaffected by the device’s control actions. Similar challenges appear for autonomous systems operating in unknown stochastic environments, or economic systems influenced by latent market factors. In such settings, designing optimal control strategies requires accounting for both the predictable evolution of the system and the stochastic nature of the surrounding environment. Effectively controlling such systems requires models that capture both the deterministic evolution of the agent’s state and the stochastic evolution of the environment. To address this challenge, in this work, we model the agent (e.g., battery energy storage system, self-driving car) with deterministic linear dynamics derived from first principles, while we model the environment (e.g., net load demand, traffic) as a linear Markov process.

Classical control theory offers elegant solutions for systems with entirely known dynamics, such as the Linear Quadratic Regulator (LQR) optimal control problem, which yields a closed-form optimal policy via Riccati equations (Anderson & Moore, 2007). On the other hand, reinforcement learning (RL) approaches have developed data-driven techniques for decision-making in unknown environments, including model-based approaches. One compelling model in this setting is that of linear Markov decision processes (linear MDPs) that leverage feature-based representations to approximate the transition kernel (Bradtke & Barto, 1996;

Francisco & Ribeiro, 2007; Sutton & Barto, 2018). The linearity of the Markov kernel coupled with the non-linearity of feature functions results in a rich but tractable model.

Yet, even such a tractable framework does not distinguish between the system dynamics and environment, viewing them as a single entity driven by the same dynamics. To this end, we propose a RL framework that combines both the LQR and linear MDP paradigms: a deterministic, discrete-time, linear time-invariant (LTI) system coupled with a stochastic environment modeled as a feature-based linear Markov process with unknown transition kernel that is unaffected by the control actions. Our objective is to design a controller that minimizes a quadratic cost over the joint system and environment states and the control actions. We first consider a nominal setting in which the LTI dynamics and the quadratic cost weights are known a priori, while only the environment’s transition kernel is unknown. We also consider extensions of the nominal setting to systems with Gaussian process noise and to LTI systems with unknown dynamics.

Our approach: By combining the structure of LQR and linear MDPs, we derive parametrized closed-form expressions for the optimal state-action value function and the corresponding optimal policy, capturing both LTI dynamics and latent stochastic effects of the environment within a unified model. This hybrid model preserves the simplicity of LQR policies while incorporating the expressive stochastic modeling of linear MDPs. Based on this structure, we develop a least-squares value iteration (LSVI) algorithm to learn the value function parameters from online data in an episodic fashion. The closed-form expression of the policy that optimizes the state-action value function makes it amenable to efficiently perform the policy update directly using the updated parameters at the end of each episode. Furthermore, because the unknown transition kernel is unaffected by the control actions, our LSVI algorithm does not require exploration in the nominal setting. For the nominal setting, we establish closed-loop stability guarantees under the learned policy and show that our LSVI achieves a regret bound of $\tilde{O}\left(T\sqrt{dL}\right)$ with high probability, where d , T , and L denote the dimension of the feature-space of the linear Markov model, the time horizon of each episode, and the number of episodes, respectively. We further extend our framework to systems with Gaussian process noise and to systems with unknown dynamics. These extensions are illustrated numerically. However, stability and convergence analysis for these settings are left for future work.

1.1 Related work

Our work lies at the intersection of optimal control and reinforcement learning, where we bridge ideas from the LQR optimal control problem and linear MDPs. Below, we review prior work that has been done in each area and highlight how our approach uniquely integrates them.

Linear Quadratic Regulator (LQR): The classical LQR problem admits closed-form optimal control policies for linear systems. Traditional methods assume full knowledge of the system dynamics and cost, enabling the computation of optimal policies via Riccati equations (Anderson & Moore, 2007). Recent work has studied the LQR problem in data-driven settings. Direct data driven approaches have been studied in (De Persis & Tesi, 2020; Dörfler et al., 2022; Celi et al., 2023), where the optimal policy is learned directly from offline data generated by the open-loop system. Indirect data-driven approaches, explored in (Aangenent et al., 2005; da Silva et al., 2018; Dean et al., 2020), first identify a model of the system dynamics from data then solve the LQR problem using the identified model. Other works have studied the LQR problem in online learning setting (Fazel et al., 2018; Mohammadi et al., 2019; Bu et al., 2019; Fatkhullin & Polyak, 2021; Bradtke et al., 1994), where the optimal policy is learned online using policy gradient methods.

Reinforcement learning with function approximation: In many reinforcement learning (RL) problems, the state or action spaces are too large (or continuous) to allow for tabular representations of value functions or policies (Kober et al., 2013; Mnih et al., 2013; Silver et al., 2016). To address this, function approximation techniques are employed to generalize from observed states and actions to unseen ones, enabling scalability and improved sample efficiency. Among the function approximation models, linear function approximation is particularly appealing due to its computational simplicity, theoretical tractability, and its ability to support efficient learning algorithms. Early approaches such as temporal difference learning, Q-learning, and least-squares temporal difference (LSTD) algorithms with linear value function approximation were explored in works like (Bradtke & Barto, 1996; Francisco & Ribeiro, 2007; Sutton & Barto, 2018).

While these methods laid important foundations, they often lacked sample efficiency guarantees and relied on heuristic exploration. Recent studies have introduced sample-efficient algorithms for linear MDPs, where the transition kernel is assumed to be a linear function of known features and unknown parameters (Yang & Wang, 2019; 2020; Jin et al., 2020). In (Jin et al., 2020), the authors developed a sample-efficient reinforcement learning algorithm for linear MDPs with a finite action space and a potentially infinite state space. Their model represents the transition kernel as a linear combination of known features with unknown probability measures, and assumes the reward function is linear in the same features with unknown parameters. In (Yang & Wang, 2020), the authors proposed a sample-efficient reinforcement learning algorithm under a linear MDP setting with possible infinite state and action space. In their framework, they assume the reward is known; further, their model introduces an additional structural assumption compared to (Jin et al., 2020), by parameterizing the transition kernel with a low-dimensional unknown matrix. This assumption reduces the learning problem to estimating this matrix, thereby significantly lowering the overall learning complexity. In our framework, we model the stochastic environment as a feature-based linear Markov Process. Similar to (Jin et al., 2020), we represent the transition kernel as a linear combination of known features with unknown probability measures. However, we assume a known quadratic cost that is independent of the features and consistent with the LQR framework, enabling efficient policy computation. This choice of the cost is more realistic and aligns with common formulations in engineering applications. Furthermore, our model avoids explicit parametric assumptions on the transition kernel made in (Yang & Wang, 2020), while allowing infinite state and action spaces. Additionally, our approach bypasses full model estimation by learning the value function directly through least-squares, benefiting from control-theoretic structure to ensure stability as well as computational and sample efficiency. Finally, our framework does not require exploration, since the environment is exogenous and is unaffected by the control inputs. Beyond linear MDPs, several works propose generalized model classes for sample-efficient RL including Bellman rank class, (Jiang et al., 2017), linear Bellman-complete classes (Munos, 2005; Zanette et al., 2020), witness rank (Sun et al., 2019), and bilinear class (Du et al., 2021).

1.2 Contributions

We list our contributions below.

- We propose a RL framework that unifies the classical LQR optimal control problem with linear MDPs. This hybrid model captures both the deterministic dynamics of physical systems and the stochastic evolution of exogenous environments. To the best of our knowledge, this integration has not been addressed in the existing literature.
- We derive a parametric form for the optimal state-action value function that decouples the agent’s dynamics from the environment’s stochasticity. This yields a closed-form policy that exhibits the simplicity of the LQR while inheriting the rich modeling capabilities linear MDPs.
- We propose a least-squares value iteration (LSVI) algorithm that learns the optimal policy by directly estimating the value function parameters. The LQR structure of our problem allows the control policy to be explicitly expressed in terms of the learned parameters, without requiring optimization of the value function at each step as in (Jin et al., 2020), thus simplifying the algorithm’s computational complexity.
- For the nominal setting, in which the linear system dynamics are known, we establish stability guarantees for the closed-loop system under the learned policy. These guarantees are given in terms of input-to-state stability, extending beyond standard sample-efficiency results in RL literature, where it is typically assumed that the reward is bounded.
- For the nominal setting, we derive a regret bound for our LSVI algorithm, showing that it achieves a rate $\tilde{\mathcal{O}}\left(T\sqrt{dL}\right)$ with high probability, where d , T , and L denote the dimension of the feature-space of the linear Markov model, the time horizon of each episodes, and the number of episodes, respectively.
- Going beyond the nominal setting, we extend our framework to systems with Gaussian process noise and to systems with unknown linear dynamics.

- We provide numerical examples to demonstrate the effectiveness of our framework in the nominal, noisy, and unknown system dynamics settings. In the nominal setting, we highlight the convergence and verify the closed-loop stability of the learned policy.

2 Problem formulation

Consider an agent obeying the discrete-time, linear, time-invariant dynamics over a finite time horizon

$$x_{t+1} = Ax_t + Bu_t, \quad t \in \{0, 1, \dots, T-1\}, \quad (1)$$

where $x_t \in \mathcal{X} = \mathbb{R}^n$ denotes the state and $u_t \in \mathcal{U} = \mathbb{R}^m$ the input with $x_0 \sim \mathcal{N}(0, \Sigma_x)$ with $\Sigma_x \succ 0$. We assume the linear dynamical system, defined via the matrix pair (A, B) , is controllable¹. We consider an environment evolving according to the discrete-time Markov Process

$$s_{t+1}|s_t \sim \mathbb{P}_t(s_{t+1}|s_t), \quad t \in \{0, \dots, T-1\}, \quad (2)$$

where $s \in \mathcal{S} \subset \mathbb{R}^p$ denotes the state of the Markov Process and $\mathbb{P}_t(s'|s)$ denotes the transition probability from state s to s' , with $s_0 \sim \mu_0$ for some distribution $\mu_0 \in \Delta(\mathcal{S})$, where $\Delta(\mathcal{S})$ denotes the set of distributions over \mathcal{S} . We assume that the matrices A and B in eq. (1) are known, while the transition probability, \mathbb{P}_t , in eq. (2) is unknown. The agent follows a control policy $\pi_t : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{U}$, where $u_t = \pi_t(x_t, s_t)$ is the action that the agent takes at state x_t and s_t at time t , for $t \geq 0$. The objective is to find an optimal control policy, $\pi = (\pi_0, \dots, \pi_T)$, that optimizes the following control task

$$\begin{aligned} & \underset{\pi}{\text{minimize}} && \mathbb{E} \left[\sum_{t=0}^T c(x_t, s_t, u_t) \right], \\ & \text{subject to} && x_{t+1} = Ax_t + Bu_t, \\ & && s_{t+1} \sim \mathbb{P}_t(s_{t+1}|s_t), \\ & && u_t = \pi_t(x_t, s_t), \end{aligned} \quad (3)$$

where $c : \mathcal{X} \times \mathcal{S} \times \mathcal{U} \rightarrow \mathbb{R}_{\geq 0}$ is the cost evaluated at x_t, s_t , and u_t for $t \geq 0$ with $u_T = \pi_T(x_T, s_T) = 0$. We restrict our search in eq. (3) to the class of deterministic policies. We show later in Section 3.1 that the optimizer of eq. (3) is indeed deterministic. We introduce the following assumptions on the transition probability in eq. (2) and the cost in eq. (3).

Assumption 2.1. (Linear Markov Process) Let $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$ be a known feature vector and $\mu_t \in \mathbb{R}^d$ a vector of d unknown signed measures over \mathcal{S} . For $s', s \in \mathcal{S}$, we have

$$\mathbb{P}_t(s'|s) = \phi(s)^\top \mu_t(s'), \quad (4)$$

We assume $\|\phi(s)\| \leq 1/\sqrt{d}$ and $\|s\| \leq \delta_s$ for all $s \in \mathcal{S}$, $\mathbb{E}[\phi(s_t)\phi(s_t)^\top] \succ 0$, and $\|\mu_t\| \leq 1$ for all t .

Assumption 2.2. (Quadratic cost) For $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $u \in \mathcal{U}$, we have

$$c(x, s, u) = \begin{bmatrix} x \\ s \\ u \end{bmatrix}^\top \underbrace{\begin{bmatrix} W & F & D \\ F^\top & M & H \\ D^\top & H^\top & R \end{bmatrix}}_C \begin{bmatrix} x \\ s \\ u \end{bmatrix}, \quad (5)$$

where $C \succeq 0$ is known and $R \succ 0$. Further, we assume the pair $(A, W^{1/2})$ is observable.

Assumption 2.1 is inspired by the linear MDP framework introduced in (Bradtke & Barto, 1996; Francisco & Ribeiro, 2007; Jin et al., 2020). However, unlike the original definition, our model assumes that the stochastic process governs only the exogenous state and is unaffected by control input. This assumption is motivated by the fact that, in our target applications, the environment is not influenced by control actions. Moreover, it

¹When the system is controllable, it implies that there exist an input sequence, u , that can drive the system from its initial state, x_0 to any final state, x_t , within finite time horizon (see (Ogata, 2010, Section 9.8)).

simplifies the expression of the optimal policy, as the optimal policy requires minimizing a quadratic function in the input u (from Assumption 2.2), rather than the nonlinear (possibly non-convex) function ϕ . We define the value function $V_t^\pi : \mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}$ as the expected cumulative cost incurred under policy π starting from state x_t and s_t at time $t \geq 0$, given by

$$V_t^\pi(x, s) \triangleq \mathbb{E} \left[\sum_{i=t}^T c(x_i, s_i, \pi_i(x_i, s_i)) \mid x_t = x, s_t = s \right].$$

Further, we define the state-action value function $Q_t^\pi : \mathcal{X} \times \mathcal{S} \times \mathcal{U} \rightarrow \mathbb{R}$ as the expected cumulative cost under policy π starting from state x_t , s_t , and action u_t at time $t \geq 0$, given by

$$Q_t^\pi(x, s, u) \triangleq c(x, s, u) + \mathbb{E} \left[\sum_{i=t+1}^T c(x_i, s_i, \pi_i(x_i, s_i)) \mid x_t = x, s_t = s, u_t = u \right].$$

To learn the optimal policy, we focus on estimating the state-action value function Q_t^π , since it directly guides policy improvement through greedy action selection. In particular, by learning an appropriate parametric approximation of the state-action value function, Q_t^π , we can infer an optimal policy without explicitly learning the transition probability measures, μ , in eq. (4). This approach leverages the structure of the system and cost, allowing us to bypass the need for full system identification and instead focus on value function approximation within the RL framework.

Remark 1. (On the knowledge of A and B in eq. (1)) *In many control applications—such as robotics and power systems—the plant dynamics (i.e., A and B) can be easily derived from first principles or can be accurately identified through standard system identification techniques prior to deployment. Our framework leverages this knowledge to focus on learning the stochastic environment component, which simplifies the computational complexity of the policy update step, and enables stability-aware control without requiring aggressive exploration. Nonetheless, we extend our framework to settings in which A and B are unknown in Section 5. We also extend it to systems with Gaussian process noise in Section 4.*

3 Main Results

We leverage the linear structures of the system in eq. (1), the transition model in eq. (4), and the quadratic structure of the cost in eq. (5) to derive a parametric expression for the state-action value function that is linear in the feature map, ϕ , along with a parametric expression for the optimal greedy policy. We introduce a least-squares value iteration algorithm to learn the parameters of the state-action value function, and therefore learn the optimal policy. We provide stability guarantees for the closed-loop system under the learned policy and a convergence analysis yielding a high-probability regret bound.

3.1 State-action value function approximation

Let the optimal value function at time t and evaluated at $x \in \mathcal{X}$ and $s \in \mathcal{S}$ under the optimal policy, π_t^* , be denoted by $V_t^*(x, s)$. Following the Bellman optimality equation, we can write the optimal state-action value function at time t and evaluated at $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $u \in \mathcal{U}$ under π_t^* as

$$Q_t^*(x, s, u) = c(x, s, u) + \mathbb{E}_{s' \sim \mathbb{P}_t(s'|s)} \{ V_{t+1}^*(Ax + Bu, s') \mid s \}.$$

The next result provides an explicit parametric form for the state-action value function Q_t .

Theorem 3.1. (Q -function representation) *Consider the dynamics in eq. (1) and the Markov Process in eq. (2). Let Assumption 2.1 and Assumption 2.2 be satisfied. Then, for any $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $u \in \mathcal{U}$, and under π_t^* for $t \geq 0$, there exists $\bar{h}_{i,t+1} \in \mathbb{R}^n$ and $\bar{q}_{i,t+1} \in \mathbb{R}$ such that*

$$Q_t^*(x, s, u) = c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \left(2(Ax + Bu)^\top \bar{h}_{i,t+1} + \bar{q}_{i,t+1} \right), \quad (6)$$

where G solves the discrete-time algebraic Riccati equation

$$G_t = A^\top G_{t+1} A + W - (A^\top G_{t+1} B + D)(R + B^\top G_{t+1} B)^{-1} (B^\top G_{t+1} A + D^\top), \quad (7)$$

with $G_T = W$.

A proof of Theorem 3.1 is in Appendix A. Several comments are in order. First, by leveraging the linearity of the system in eq. (1) and the Markov process in eq. (4), along with the quadratic structure of the cost in eq. (5), the state-action value function in eq. (6) exhibits a structure that decouples the linear system state x and action u from the exogenous state s . Second, the derived expression of the state-action value function in eq. (6) is linear in the feature map ϕ and the weight parameters $\bar{h}_{i,t}$ and $\bar{q}_{i,t}$. Third, the weight parameters $\bar{h}_{i,t}$ and $\bar{q}_{i,t}$ depend on the unknown transition probability $\mathbb{P}(\cdot|s)$ in eq. (4), and therefore, learning the state-action value function boils down to learning these weights, thereby bypassing the need to explicitly learn the probability measures, μ , in eq. (4). The optimal policy is found by minimizing the Q function over the input u . Since by Theorem 3.1, Q is quadratic in u , this optimal policy can be found in closed form, as shown in the following corollary, which expresses the optimal policy in terms of feedback gains and weight parameters.

Corollary 3.2. (Optimal policy representation) For any $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $t \in \{0, 1, \dots, T-1\}$

$$u_t^*(x, s) = \pi_t^*(x, s) = K_{x,t}x + K_{s,t}s + K_{h,t} \sum_{i=1}^d \phi_i(s) \bar{h}_{i,t+1},$$

where $\bar{h}_{i,t+1}$ is as in Theorem 3.1 and

$$\begin{aligned} K_{x,t} &= -(R + B^\top G_{t+1} B)^{-1} (B^\top G_{t+1} A + D^\top), \\ K_{s,t} &= -(R + B^\top G_{t+1} B)^{-1} H^\top, \\ K_{h,t} &= -(R + B^\top G_{t+1} B)^{-1} B^\top, \end{aligned} \quad (8)$$

and G_{t+1} satisfies eq. (7).

Algorithm 1 Least-Squares Value Iteration

- 1: Given: L, R_θ, λ
 - 2: **for** episode $\ell = 1, \dots, L$ **do**
 - 3: $x_0^\ell \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_x)$ with $\Sigma_x \succ 0$
 - 4: $s_0^\ell \stackrel{\text{i.i.d.}}{\sim} \mu_0$ such that $\mathbb{E}[\phi(s_0)\phi(s_0)^\top] \succ 0$
 - 5: **for** step $t = T-1, \dots, 0$ **do**
 - 6: $\Lambda_t^\ell \leftarrow \sum_{i=1}^{\ell-1} Y(x_t^i, u_t^i)^\top \phi(s_t^i) \phi(s_t^i)^\top Y(x_t^i, u_t^i) + \lambda I_{dn+d}$
 - 7: $\theta_{t+1}^\ell \leftarrow (\Lambda_t^\ell)^{-1} \sum_{i=1}^{\ell-1} Y(x_t^i, u_t^i)^\top \phi(s_t^i) \epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i)$
 - 8: **if** $\|\theta_{t+1}^\ell\| > R_\theta$ **then**
 - 9: $\theta_{t+1}^\ell \leftarrow \frac{R_\theta}{\|\theta_{t+1}^\ell\|} \theta_{t+1}^\ell$
 - 10: **end if**
 - 11: **end for**
 - 12: **for** step $t = 0, \dots, T-1$ **do**
 - 13: $u_t^\ell \leftarrow K_{x,t}x_t^\ell + K_{s,t}s_t^\ell + K_{h,t}(\phi(s_t^\ell)^\top \otimes Z)\theta_{t+1}^\ell$
 - 14: Take action u_t^ℓ
 - 15: Observe x_{t+1}^ℓ and s_{t+1}^ℓ
 - 16: **end for**
 - 17: **end for**
-

3.2 Learning weight parameters of the value function via least-squares value iteration

In this subsection, we learn the weight parameters, \bar{h} and \bar{q} , that parameterize the state-action value function in Theorem 3.1. To this aim, we propose a least-squares value iteration algorithm (Algorithm 1) that is inspired by (Jin et al., 2020). Before we lay out the steps of our algorithm, we introduce the following notations. At each time step, t , we concatenate the parameters $\bar{h}_{i,t}$ and $\bar{q}_{i,t}$ for $i \in \{1, \dots, d\}$ as

$$\theta_t = [\theta_{1,t}^\top \ \cdots \ \theta_{d,t}^\top]^\top, \quad \text{where} \quad \theta_{i,t} = [\bar{h}_{i,t}^\top \ \bar{q}_{i,t}^\top]^\top. \quad (9)$$

Using the notation in eq. (9), we rewrite the Q -function in Theorem 3.1 and the policy in Corollary 3.2 as

$$\begin{aligned} Q_t(x, s, u) &= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \phi(s)^\top Y(x, u) \theta_{t+1}, \\ u_t(x, s) &= K_{x,t}x + K_{s,t}s + K_{h,t} \left(\phi(s)^\top \otimes Z \right) \theta_{t+1}, \end{aligned} \quad (10)$$

where $Y(x, u) = I_d \otimes [2(Ax + Bu)^\top, 1]$ and $Z = [I_n, 0_{n \times 1}]$. Now we lay out the steps of our least-squares value iteration algorithm (Alg. 1). Our algorithm consists of an outer loop over L episodes, where each episode consists of two loops: 1) backward-in-time weight update loop (lines 5-11) and 2) forward roll-out and data collection loop (lines 12-16). During the first pass of episode ℓ (lines 5-11), we treat the data collected in the previous $\ell - 1$ episodes as a fixed dataset

$$\mathcal{D}_{\ell-1} := \{(x_t^i, s_t^i, u_t^i, x_{t+1}^i, s_{t+1}^i) : i < \ell, 0 \leq t < T\}. \quad (11)$$

At each time step t , θ minimizes a regularized least-squares loss—the squared error between the parametric state-action value function in eq. (10) and the Bellman target (immediate cost plus the estimated value of the next state). Solving this problem on past trajectory data yields an accurate value-function approximation and enables closed-form greedy policy updates without estimating the transition probabilities. The regularized least-squares regression is stated as (see Appendix B for details)

$$\theta_{t+1}^\ell = \arg \min_{\theta \in \mathbb{R}^{d(n+1)}} \sum_{i=1}^{\ell-1} \left(\phi(s_t^i)^\top Y(x_t^i, u_t^i) \theta - \epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i) \right)^2 + \lambda \|\theta\|^2,$$

where

$$\epsilon_{t+1}^\ell(x, s) = 2(x)^\top h_{t+1}^\ell(s) + q_{t+1}^\ell(s), \quad (12)$$

$$h_{t+1}^\ell(s) = (A^\top + K_{x,t}^\top B^\top) (\phi(s)^\top \otimes Z) \theta_{t+2}^\ell + (F + K_{x,t}^\top H^\top) s, \quad (13)$$

$$\begin{aligned} q_{t+1}^\ell(s) &= \left(\phi(s)^\top \otimes \bar{Z} \right) \theta_{t+2}^\ell + s^\top (M + HK_{s,t}) s + \theta_{t+2}^{\ell-1} \left(\phi(s) \otimes Z^\top \right) BK_{h,t} \phi(s)^\top \otimes Z \theta_{t+2}^\ell \\ &\quad + 2s^\top HK_{h,t} \left(\phi(s)^\top \otimes Z \right) \theta_{t+2}^\ell, \end{aligned} \quad (14)$$

with $\bar{Z} = [0_{1 \times n}, 1]$. Unlike prior work (e.g., (Jin et al., 2020)), we leverage the structure of our model to derive a closed-form expression for the Bellman target in terms of previously learned parameters, thereby avoiding an inner optimization over the action space at each time step (often required in discrete action space settings). In fact, $\epsilon_{t+1}^\ell(x, s)$ is obtained directly from this closed-form Bellman target (see Appendix B). The closed-form parameter update is given by

$$\begin{aligned} \Lambda_t^\ell &= \sum_{i=1}^{\ell-1} Y(x_t^i, u_t^i)^\top \phi(s_t^i) \phi(s_t^i)^\top Y(x_t^i, u_t^i) + \lambda I_{d(n+1)}, \\ \theta_{t+1}^\ell &= (\Lambda_t^\ell)^{-1} \sum_{i=1}^{\ell-1} Y(x_t^i, u_t^i)^\top \phi(s_t^i) \epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i), \end{aligned} \quad (15)$$

which recover lines 6 and 7 of Alg. 1. For $\ell = 1$, we set $\theta_{t+1}^\ell = 0$ and $\Lambda_t^\ell = \lambda I_{d(n+1)}$ for $t \in \{0, \dots, T-1\}$. The regularizer term $\lambda I_{d(n+1)}$ ensures numerical stability, the projection step in lines 8-10 makes sure that

the norm of the learned parameters is uniformly bounded for $t \in \{0, \dots, T-1\}$ and $\ell \in \{1, \dots, L\}$. In the second pass (lines 12–16) the newly computed parameters θ_{t+1}^ℓ are plugged into the policy eq. (10),

$$u_t^\ell(x_t^\ell, s_t^\ell) = K_{x,t}x_t^\ell + K_{s,t}s_t^\ell + K_{h,t}(\phi(s_t^\ell)^\top \otimes Z)\theta_{t+1}^\ell,$$

generating a new trajectory $(\{x_t^\ell, s_t^\ell, u_t^\ell\}_{t=0}^T)$. These samples are appended to the collected data eq. (11), and will be used in the next episode’s backward update. Notice that Alg. 1 does not require exploration as in (Jin et al., 2020), which we discuss in the following remark.

Remark 2. (Role of exploration) *In classical reinforcement learning, exploration (e.g., using ε -greedy or optimism-based methods) is necessary to sufficiently explore the environment and estimate unknown transition dynamics. However, our framework does not require exploration. This is because the stochastic component of the environment is modeled as an exogenous Markov process that evolves independently of the control inputs (see Assumption 2.1), and the system dynamics, A and B , are known. Our algorithm estimates the value function parameters, θ , via a least-squares procedure using observed trajectories without the need to infer the transition probabilities explicitly. Thus, the optimal policy can be computed in closed form by minimizing a known quadratic function of the input.*

Remark 3. (Choice of R_θ) *The projection radius R_θ in Alg. 1 ensures that the learned parameters at each episode and time step remain within a ball of radius R_θ , ensuring numerical stability. Moreover, it plays a crucial role in the theoretical analysis (i.e., stability and regret bound). In practice, R_θ should be chosen large enough to contain the true parameters, θ^* , but not excessively large to keep the constants in the stability and regret bounds moderate. In Appendix C, we derive an upper bound on $\|\theta^*\|$; if R_θ is larger than this bound, then the ball of radius R_θ is guaranteed to contain θ^* . In particular, we show that θ^* is contained in this ball if $R_\theta \geq c_\theta\sqrt{d}$, where $c_\theta > 0$ depends on known problem parameters, e.g., the system matrices, cost weights, the feature map, and the bound on s .*

3.3 Input-to-State Stability

It is critical to ensure that the learned policy stabilizes the closed-loop system in each episode, particularly in settings where the environment evolves independently of the control actions and safety is a concern. To this end, we establish an input-to-state stability (ISS) bound for the system under the learned policy, which we present in the following result.

Theorem 3.3. (Input-to-state stability) *Consider system eq. (1), let u be the output of Algorithm 1 at episode ℓ . Let $\|\theta_t^\ell\| \leq R_\theta$, $\|K_{s,t}\| \leq \bar{K}_s$, and $\|K_{h,t}\| \leq \bar{K}_h$ for $t \in \{0, \dots, T-1\}$. Let x_0^ℓ be the initial state in episode ℓ . Then, under Assumptions 2.1 and 2.2,*

$$\|x_t^\ell\| \leq \alpha\rho^t\|x_0^\ell\| + \frac{\alpha\|B\|}{1-\rho} \left(\bar{K}_s\delta_s + \frac{\bar{K}_h R_\theta}{\sqrt{d}} \right), \quad (16)$$

for $t \in \{0, \dots, T-1\}$, where $\alpha > 0$ and $0 < \rho < 1$ are constants.

A proof of Theorem 3.3 is deferred to Appendix D. Several comments are in order. First, Theorem 3.3 implies that the state trajectory at each episode $\ell \in \{1, \dots, L\}$ remains bounded in terms of the initial condition, the system dynamics, the control gains in Corollary 3.2, and R_θ . Second, the first term on the right-hand side of eq. (16), which depends on the initial state, decays exponentially with time, while the second term is independent of time and the number of episodes. This latter term depends on the system matrices, feedback gains, the bound on s , R_θ from Algorithm 1, and the dimension of the feature map, d . This result leverages the known system dynamics and the structure of the policy, extending traditional stability notions in control to learning-based policies in partially known settings.

3.4 Regret analysis

We define the regret $\mathcal{R}(L)$ as the difference between the total cost incurred by the learned policy and that of the optimal policy over L episodes. Mathematically, for L episodes, the regret is defined as

$$\mathcal{R}(L) = \sum_{\ell=1}^L (V_0^\ell(x_0^\ell, s_0^\ell) - V_0^*(x_0^\ell, s_0^\ell)). \quad (17)$$

where $V_0^\ell(x_0^\ell, s_0^\ell)$ denotes the value evaluated at the initial states x_0^ℓ and s_0^ℓ under the policy learned at episode ℓ , and $V_0^*(x_0^\ell, s_0^\ell)$ is the value of the optimal policy evaluated at the initial states x_0^ℓ and s_0^ℓ . We derive a bound on the regret in the following result.

Theorem 3.4. (Regret bound) *Let Assumptions 2.1 and 2.2 be satisfied. Let $\|Y(x_t^\ell, u_t^\ell)^\top \phi(s_t^\ell)\| \leq \delta_\psi$, for $t \in \{0, \dots, T\}$ and $\ell \in \{1, \dots, L\}$. Let $\beta = \log\left(1 + \frac{L\delta_\psi^2}{\lambda}\right)$ with $\lambda > 0$. Let $\delta \in [0, 1/3]$. Then, with probability at least $1 - 3\delta$*

$$\mathcal{R}(L) \leq \sigma \sqrt{2TL \log(1/\delta)} + \delta_\psi T \left(\frac{1}{\sqrt{\lambda}} + \frac{4\sqrt{L}}{\sqrt{\gamma}} \right) \left(\sigma \sqrt{2dn\beta + 2 \log\left(\frac{1}{\delta}\right)} + (R_\theta + 2\delta_v)\sqrt{\lambda} \right)$$

where $\sigma > 0$, $\gamma > 0$ and $\delta_v > 0$ are constants that do not depend on L and T , and do not scale with d . Further, δ_ψ scales with $\mathcal{O}(1/\sqrt{d})$ and R_θ scales with $\mathcal{O}(\sqrt{d})$.

A proof of Theorem 3.4 is presented in Appendix E. Theorem 3.4 provides a probabilistic upper bound on the cumulative regret. Several comments are in order. First, the leading term of the bound scales as $\mathcal{O}\left(T\sqrt{dL \log(L)}\right)$ or $\tilde{\mathcal{O}}\left(T\sqrt{dL}\right)$, which matches, in terms of the number of episodes, the rate reported in (Jin et al., 2020). Second, our bound grows linearly in T and \sqrt{d} , in contrast to the T^2 and $d\sqrt{d}$ factors in (Jin et al., 2020), respectively. Third, the constants σ and δ_v are independent of L and T , and they do not scale with d and they depend only on the system matrices in eq. (1), the cost weight matrices in eq. (5), the bound on the state, x_t , in eq. (16), and the bound on the exogenous state s_t in eq. (4). In fact, they arise from the uniform upper bound on the value function, V_t (see Appendix E for the explicit formulas). Fourth, the constant γ satisfies $\mathbb{E}[Y(x_0, u_0)^\top \phi(s_0) \phi(s_0)^\top Y(x_0, u_0)] \succeq \gamma I_{d(n+1)}$, which holds because the initial states, x_0 and s_0 are drawn independently in each episode. Finally, the bound in Theorem 3.4 suggests, when the initial states, x_0 and s_0 , are fixed for all episodes, Alg. 1 can learn an ε -optimal policy, π , that satisfies $V_0^\pi(x_0, s_0) - V_0^*(x_0, s_0) \leq \varepsilon$ after $L = \tilde{\mathcal{O}}\left(\frac{dT^2}{\varepsilon^2}\right)$ episodes.

4 Extension to Systems with Process Noise

In this section, we add Gaussian process noise to the system in eq. (1). In particular, we consider the following discrete-time, linear, time-invariant dynamics over a finite time horizon

$$x_{t+1} = Ax_t + Bu_t + w_t, \quad t \in \{0, 1, \dots, T-1\}, \quad (18)$$

where $x_t \in \mathcal{X} = \mathbb{R}^n$ denotes the state, $u_t \in \mathcal{U} = \mathbb{R}^m$ the input with $x_0 \sim \mathcal{N}(0, \Sigma_x)$ with $\Sigma_x \succ 0$, and $w_t \in \mathbb{R}^n$ the process noise, where $w \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_w)$ and $\Sigma_w \succ 0$.

4.1 State-action value function approximation

Following the Bellman optimality equation, we can write the optimal state-action value function at time t and evaluated at $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $u \in \mathcal{U}$ under π_t^* with $w_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_w)$ as

$$Q_t^*(x, s, u) = c(x, s, u) + \mathbb{E}_w \left\{ \mathbb{E}_{s' \sim \mathbb{P}_t(s'|s)} \left\{ V_{t+1}^*(Ax + Bu + w, s') \mid s \right\} \right\}.$$

The next result provides an explicit parametric form for the state-action value function Q_t and V_t .

Theorem 4.1. (Value-function representation) *Consider the dynamics in eq. (18) and the Markov Process in eq. (2). Let Assumption 2.1 and Assumption 2.2 be satisfied. Then, for any $x \in \mathcal{X}$, $s \in \mathcal{S}$, and $u \in \mathcal{U}$, and under π_t^* for $t \geq 0$, there exists $\bar{h}_{i,t+1} \in \mathbb{R}^n$ and $\bar{q}_{i,t+1} \in \mathbb{R}$ such that*

$$Q_t^*(x, s, u) = c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \left(2(Ax + Bu)^\top \bar{h}_{i,t+1} + \bar{q}_{i,t+1} \right) + \sum_{i=t+1}^T \text{tr}[G_i \Sigma_w]. \quad (19)$$

Further,

$$V_t^*(x, s) = x^\top G_t x + 2h_t^\top(s)x + q_t(s) + \sum_{i=t+1}^T \text{tr}[G_i \Sigma_w], \quad (20)$$

where $G_t \in \mathbb{R}^{n \times n} \succ 0$, $h_t(\cdot) \in \mathbb{R}^n$, and $q_t(\cdot) \in \mathbb{R}$ are as in Theorem 3.1.

Proof. The proof follows in a similar manner as the proof of Theorem 3.1 and noting that $\mathbb{E}_w [w_t] = 0$ and $\mathbb{E}_w [w_t w_t^\top] = \Sigma_w$, for $t \in \{0, 1, \dots, T-1\}$. \square

4.2 Learning weight parameters of the value function via least-squares value iteration

In this subsection, we extend the least-squares value iteration procedure in Section 3.2 to the setting with process noise in eq. (18). The key observation is that the parametric structure of the optimal Q -function remains the same, up to an additive constant term. In particular, from Theorem 4.1, for any $x \in \mathcal{X}$, $s \in \mathcal{S}$, $u \in \mathcal{U}$ and $t \in \{0, \dots, T-1\}$, we can write

$$\begin{aligned} Q_t(x, s, u) &= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \phi(s)^\top Y(x, u) \theta_{t+1} + \sum_{i=t+1}^T \text{tr}[G_i \Sigma_w], \\ u_t(x, s) &= K_{x,t} x + K_{s,t} s + K_{h,t} \left(\phi(s)^\top \otimes Z \right) \theta_{t+1}, \end{aligned} \quad (21)$$

where θ_{t+1} and $Y(x, u) = I_d \otimes [2(Ax + Bu)^\top, 1]$ are the same as in Section 3.2. Notice that the additional term $\text{tr}(G_{t+1} \Sigma_w)$ does not depend on the input, and therefore does not affect the minimization over u . We follow the same steps as Algorithm 1. For episode ℓ , as in Section 3.2, we treat the data collected in the previous $\ell-1$ episodes as a fixed dataset $\mathcal{D}^{\ell-1}$. At each time step t , the LSVI update fits the parametrized Q -function in eq. (21) to a Bellman target. Thus, θ is obtained by solving the following regularized least-squares regression problem (see Appendix F for details):

$$\begin{aligned} \theta_{t+1}^\ell &= \arg \min_{\theta} \sum_{j=1}^{\ell-1} \left((Ax_t^j + Bu_t^j)^\top G_{t+1} (Ax_t^j + Bu_t^j) + \phi(s_t^j)^\top Y(x_t^j, u_t^j) \theta + \sum_{i=t+1}^T \text{tr}[G_i \Sigma_w] \right. \\ &\quad \left. - (x_{t+1}^j)^\top G_{t+1} x_{t+1}^j - 2 \left(h_{t+1}^\ell(s_{t+1}^j) \right)^\top x_{t+1}^j - q_{t+1}^\ell(s_{t+1}^j) - \sum_{i=t+2}^T \text{tr}[G_i \Sigma_w] \right)^2 + \lambda \|\theta\|_2^2, \end{aligned} \quad (22)$$

where $h_{t+1}^\ell(\cdot)$ and $q_{t+1}^\ell(\cdot)$ are defined as in eq. (13) and eq. (14), respectively. The closed-form parameter update is therefore given by

$$\begin{aligned} \theta_{t+1}^\ell &= \Lambda_t^{-1} \sum_{j=1}^{\ell-1} Y^\top(x^j, u^j) \phi(s^j) \epsilon^\ell(x_t^i, x_{t+1}^i, s_{t+1}^i, u_t^i), \\ \Lambda_t &= \sum_{j=1}^{\ell-1} Y(x^j, u^j)^\top \phi(s^j) \phi(s^j)^\top Y(x^j, u^j) + \lambda I_{d(n+1)}, \\ \epsilon^\ell(x_t^i, x_{t+1}^i, s_{t+1}^i, u_t^i) &= (x_{t+1}^i)^\top G_{t+1} x_{t+1}^i + 2 \left(h_{t+1}^\ell(s_{t+1}^i) \right)^\top x_{t+1}^i + q_{t+1}^\ell(s_{t+1}^i) \\ &\quad - (Ax_t^i + Bu_t^i)^\top G_{t+1} (Ax_t^i + Bu_t^i) - \text{tr}(G_{t+1} \Sigma_w). \end{aligned} \quad (23)$$

Thus, the only modification to Algorithm 1 is in Step 7, where we compute ϵ_{t+1}^ℓ as in eq. (23).

5 Extension to Unknown System Dynamics

In this section, we relax the assumption that A and B in eq. (1) are known. We begin by re-writing the Q -function in eq. (6) as

$$Q_t^*(x, s, u) = c(x, s, u) + \begin{bmatrix} x \\ u \end{bmatrix}^\top \underbrace{\begin{bmatrix} A^\top G_{t+1} A & A^\top G_{t+1} B \\ B^\top G_{t+1} A & B^\top G_{t+1} B \end{bmatrix}}_{P_{t+1}} \underbrace{\begin{bmatrix} x \\ u \end{bmatrix}}_z + 2 \sum_{i=1}^d \phi_i(s) \begin{bmatrix} x \\ u \end{bmatrix}^\top \underbrace{\begin{bmatrix} A^\top \\ B^\top \end{bmatrix}}_{\tilde{h}_{i,t+1}} \bar{h}_{i,t+1} + \sum_{i=1}^d \phi_i(s) \bar{q}_{i,t+1}. \quad (24)$$

Equation (24) preserves the same structural form as eq. (6). The difference is that in eq. (6) we assumed A and B were known, so only \bar{h}_{t+1} and \bar{q}_{t+1} were treated as unknowns. In contrast, when the dynamics are unknown, the matrices A and B become implicitly embedded inside the quantities P_{t+1} and \tilde{h}_{t+1} . Thus, the roles previously played by \bar{h}_{t+1} and \bar{q}_{t+1} in eq. (6) are now taken over by the enlarged set of parameters, P_{t+1} , \tilde{h}_{t+1} , and \bar{q}_{t+1} , all of which must be learned from data. Similarly, we re-write the policy in Corollary 3.2 as

$$u_t^*(x, s) = K_{x,t} x + K_{s,t} s + \tilde{K}_{h,t} \sum_{i=1}^d \phi_i(s) \tilde{h}_{i,21,t+1}. \quad (25)$$

where $\tilde{h}_{i,21} \in \mathbb{R}^m$ is the $(2, 1)$ -block of $\tilde{h}_i \in \mathbb{R}^{n+m}$ in eq. (24), and

$$\begin{aligned} K_{x,t} &= -(R + P_{22,t+1})^{-1} (D^\top + P_{21,t+1}), \\ K_{s,t} &= -(R + P_{22,t+1})^{-1} H^\top, \\ \tilde{K}_{h,t} &= -(R + P_{22,t+1})^{-1}, \end{aligned} \quad (26)$$

where $P_{21} \in \mathbb{R}^{m \times n}$ and $P_{22} \in \mathbb{R}^{m \times m}$ are the $(2, 1)$ -block and the $(2, 2)$ -block of the matrix P in eq. (24), respectively. The details of learning the parameters, P_t , \tilde{h}_t , and \bar{q}_t , for $t \in \{0, \dots, T\}$, together with the corresponding least-squares value iteration algorithm extending Algorithm 1 to the case where A and B are unknown, are presented in Appendix G with resulting method summarized in Algorithm 2.

Remark 4. (The need for exploration) In contrast to Algorithm 1, a key aspect of Algorithm 2 is that the input u_t in Line 22 involves an additional excitation term. This term is necessary to sufficiently excite the system so that the collected data are informative enough to learn the unknown system matrices A and B .

Remark 5. (Stability and convergence) Unlike Algorithm 1, for Algorithm 2 neither closed-loop stability nor convergence is guaranteed in general, especially when the underlying system is unstable. Since the system matrices A and B are unknown and must be learned online, the exploratory inputs and estimation errors may affect both stabilization and parameter convergence during the learning process. Therefore, further analysis is needed to establish conditions under which Algorithm 2 guarantees stability and convergence.

6 Illustration of Results

6.1 Known System Dynamics

We consider a discrete-time, linear, time-invariant system

$$x_{t+1} = \underbrace{\begin{bmatrix} 1.8 & 1.2 \\ 0 & 1.19 \end{bmatrix}}_A x_t + \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_B u_t, \quad (27)$$

and a stochastic environment evolving according to a feature-based linear Markov process with

$$s_{t+1} | s_t \sim \underbrace{\begin{bmatrix} f_1(s_t) & f_2(s_t) \\ f_1(s_t) + f_1(s_t) & f_1(s_t) + f_1(s_t) \end{bmatrix}}_{\phi(s_t)^T} \underbrace{\begin{bmatrix} \mathcal{N}(s_{t+1}; 7, 1) \\ \mathcal{N}(s_{t+1}; -1, 1.5) \end{bmatrix}}_{\mu_{t+1}(s_{t+1})}, \quad (28)$$

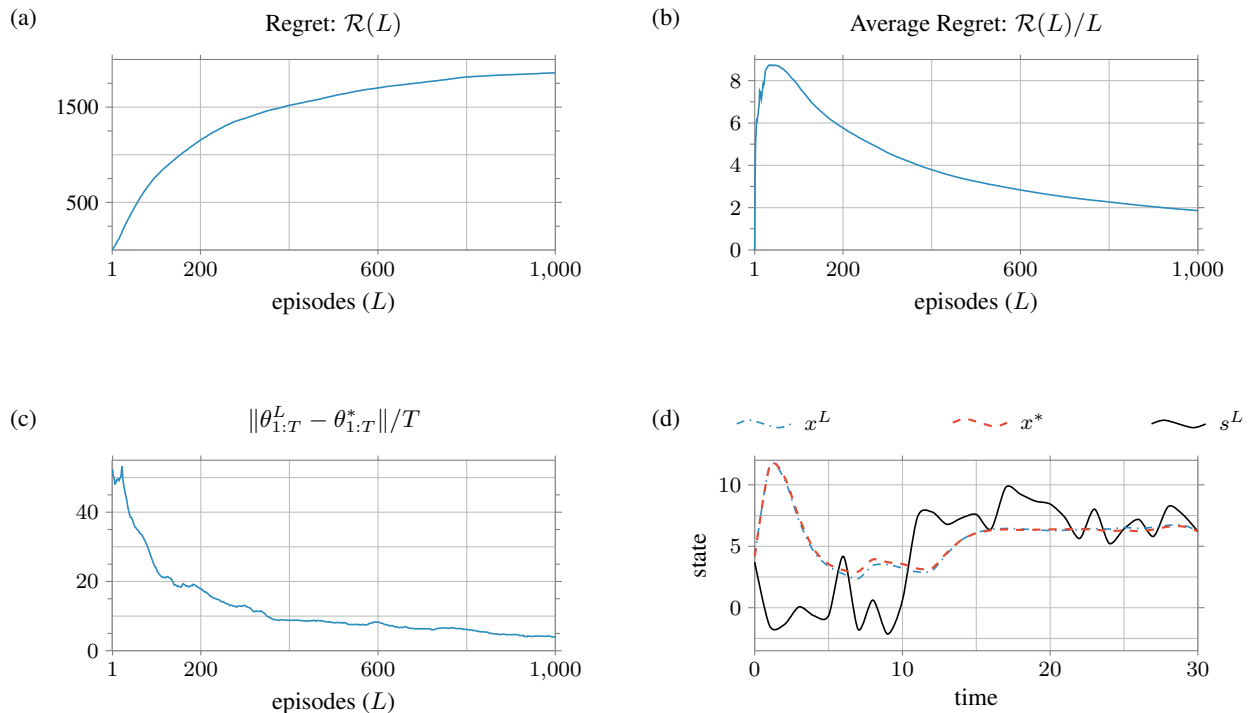


Figure 1: This figure shows the numerical results for the setting in Section 6.1. Panel (a) shows the regret as a function of L , scaling as $\tilde{O}(\sqrt{L})$ in line with Theorem 3.4. Panel (b) shows the average regret as a function of L , and we observe that it converges as L increases. Panel (c) shows the norm of the estimation error between the learned and the true parameters, averaged over the episode horizon T , as a function of L . We observe that the estimation error decreases with L , indicating that the learned policy converges to the optimal one. Panel (d) shows the state trajectory generated by the system in eq. (27) under the learned policy at episode $L = 1000$ (dot-dashed blue line) and under the optimal policy (dashed red line). It also shows the exogenous state trajectory generated by the linear Markov process in eq. (28) (solid black line). We observe that the trajectory under the learned policy closely matches that of the optimal policy, and both track the mean of the exogenous state.

where $f_1(s_t) = \exp\left(\frac{-(s_t - \nu_1)^2}{2\rho_1^2}\right)$ and $f_2(s_t) = \exp\left(\frac{-(s_t - \nu_2)^2}{2\rho_2^2}\right)$, where $\nu_1 = 7$, $\nu_2 = -1$, $\rho_1 = 5$, and $\rho_2 = 3$, with $\|s_t\| \leq \delta_s = 15$ for all t . We define the cost function to capture the tracking error between the first state of eq. (27) and the exogenous state, s , and is expressed as

$$c(x, s, u) = (Cx - s)^\top M (Cx - s) + u^\top Ru = x^\top \underbrace{C^\top MC}_W x + s^\top Ms + u^\top Ru - 2s \underbrace{MC}_{F^\top} x, \quad (29)$$

where $C = [1, 0]$, $M = 1$, and $R = 1$. First, we use the matrices A and B , along with the cost weight matrices, to compute the feedback gains in Corollary 3.2. Then, we apply Algorithm 1 to learn the parameters, θ , using $L = 1000$ episodes, each with horizon $T = 30$. We set $\lambda = 2$ and $R_\theta = 500$ (see Remark 3). At each episode, we sample $x_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 4)$ and $s_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 4)$, which are independent of each other. Using knowledge of the true distributions in eq. (28), we compute the true parameters, θ_t^* for $t \in \{1, \dots, T\}$, via the results in Appendix C.1, which we then use to compute the true optimal policy π_t^* as in Corollary 3.2. Finally, we use the true parameters, we compute the regret in eq. (17). We present our numerical results in Figure 1. The regret $\mathcal{R}(L)$ as a function of the number of episodes L is shown in Fig. 1(a). We observe that the regret scales as $\tilde{O}(\sqrt{L})$, which is consistent with our results in Theorem 3.4. The average regret, $\mathcal{R}(L)/L$ as a function of L is shown in Figure 1(b) and is observed to converge as L increases, indicating convergence of our algorithm. Fig 1(c) presents the norm of the estimation error between the learned and the true parameters, averaged

over the episode horizon T , as a function of L . This is expressed as $\|\theta_{1:T}^L - \theta_{1:T}^*\|/T$, where $\theta_{1:T}^L \in \mathbb{R}^{d(n+1) \times T}$ is a matrix whose columns corresponds to the parameters θ_t^L at each time step t , and similarly for $\theta_{1:T}^*$. We observe that the estimation error decreases with L , indicating that the learned policy gradually converges to the optimal one. Finally, we apply both the learned policy at episode $L = 1000$ and the optimal policy to the system in eq. (27), and compare their corresponding closed-loop state trajectories, as shown in Figure 1(d), alongside the trajectory of the exogenous state s . We observe that the trajectory under the learned policy closely matches that of the optimal policy, and both effectively track the mean of the exogenous state.

6.2 Known System Dynamics with Process Noise

In this subsection, we consider the same setting as in Section 6.1, except that the system in eq. (27) is affected by additive process noise. In particular, we consider

$$x_{t+1} = Ax_t + Bu_t + w_t, \quad t \in \{0, 1, \dots, T-1\}, \quad (30)$$

where $w_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_w)$ with $\Sigma_w \succ 0$. We use the same system matrices A and B as in eq. (27), the same stochastic environment in eq. (28), and the same quadratic tracking cost in eq. (29). We set $\Sigma_w = I_2$. We first use the known matrices A and B , together with the cost weight matrices, to compute the feedback gains in Corollary 3.2. Then, we implement Algorithm 1 with the modified Bellman target in eq. (23) to learn the parameters θ . We use $L = 1000$ episodes, each with horizon $T = 30$, and set $\lambda = 2$ and $R_\theta = 500$. At each episode, we sample $x_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 25)$ and $s_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 4)$, independently of each other. Since the additive term $\sum_{i=t+1}^T \text{tr}[G_i \Sigma_w]$ in Theorem 4.1 does not depend on the control input, the optimal policy retains the same closed-form structure as in Corollary 3.2, while θ is learned as in eq. (23). Using knowledge of the true distributions in eq. (28), we compute the true parameters, θ_t^* for $t \in \{1, \dots, T\}$, via the results in Appendix C.1, which we then use to compute the true optimal policy π_t^* as in Corollary 3.2. Finally, we use the true parameters, we compute the regret in eq. (17).

We present our numerical results in Figure 2. The results show that Algorithm 1 continues to learn an effective control policy in the presence of process noise. In particular, the regret $\mathcal{R}(L)$ shown in Figure 2(a) grows sublinearly with the number of episodes L . The average regret shown in Figure 2(b) decreases as L increases. Figure 2(c) presents the norm of the estimation error between the learned and the true parameters, averaged over the episode horizon T , as a function of L . This is expressed as $\|\theta_{1:T}^L - \theta_{1:T}^*\|/T$, where $\theta_{1:T}^L \in \mathbb{R}^{d(n+1) \times T}$ is a matrix whose columns corresponds to the parameters θ_t^L at each time step t , and similarly for $\theta_{1:T}^*$. We observe that the estimation error decreases with L , indicating that the learned policy gradually converges to the optimal one. Finally, we apply both the learned policy at episode $L = 1000$ and the optimal policy to the system in eq. (27), and compare their corresponding closed-loop state trajectories, as shown in Figure 2(d), alongside the trajectory of the exogenous state s . We observe that the trajectory under the learned policy closely matches that of the optimal policy, and both effectively track the mean of the exogenous state.

6.3 Unknown System Dynamics

In this subsection, we consider the same setting as in Section 6.1, except that the system matrices A and B in eq. (27) are assumed to be unknown. Consequently, we need to learn the parameters P_t , \tilde{h}_t , and \tilde{q}_t in eq. (24) jointly from data. These parameters can be combined into a single vector θ , as described in eq. (137). We use the following stable linear system (see Remark 5)

$$x_{t+1} = \underbrace{\begin{bmatrix} 0.4 & 1.2 \\ 0 & 0.19 \end{bmatrix}}_A x_t + \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_B u_t, \quad (31)$$

and the same stochastic environment and quadratic tracking cost as in eq. (28) and eq. (29), respectively. We implement Algorithm 2 with $L = 1500$ episodes and horizon $T = 30$. We set the regularization parameter to $\lambda = 2$ and the projection radius to $R_\theta = 800$. At each episode, we sample $x_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 25)$ and $s_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 25)$, independent of each other. We use an exploration term $\eta_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 25)$ for $t \in \{0, \dots, T\}$. Similarly as in

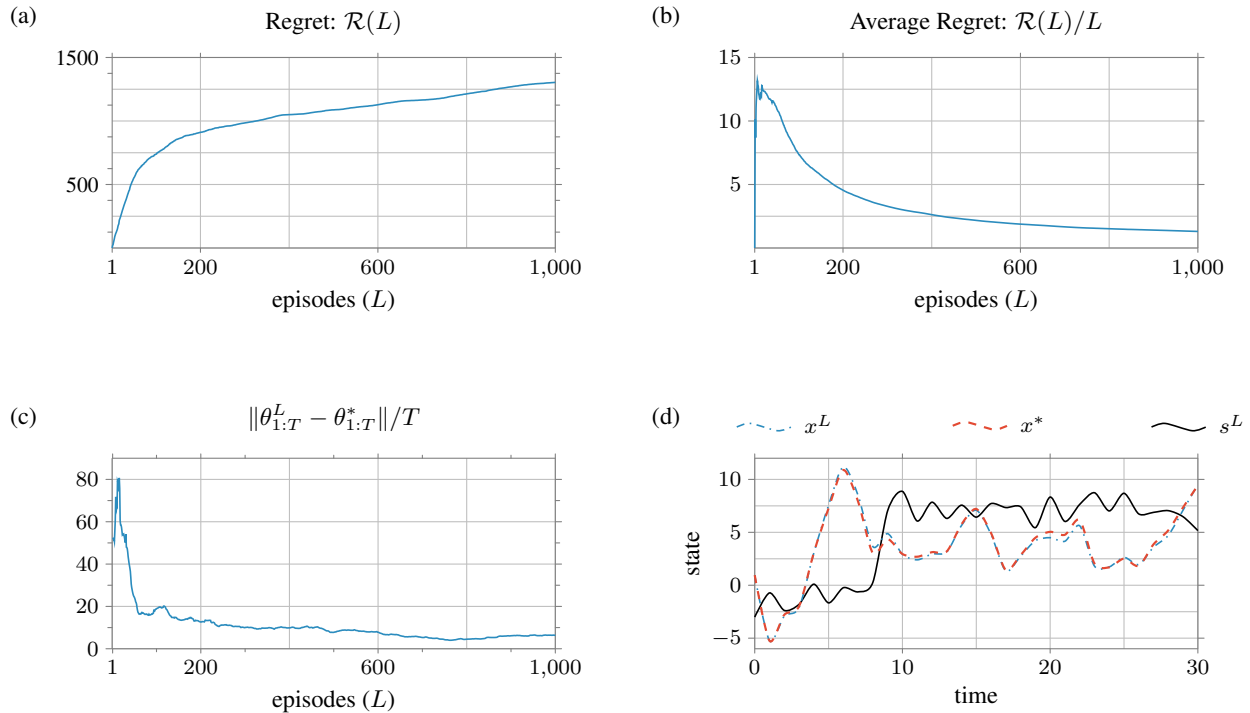


Figure 2: This figure shows the numerical results for the setting in Section 6.2. Panel (a) shows the regret as a function of L , scaling sublinearly with L . Panel (b) shows the average regret as a function of L , and we observe that it converges as L increases. Panel (c) shows the norm of the estimation error between the learned and the true parameters, averaged over the episode horizon T , as a function of L . We observe that the estimation error decreases with L , indicating that the learned policy converges to the optimal one. Panel (d) shows the state trajectory generated by the system in eq. (27) under the learned policy at episode $L = 1000$ (dot-dashed blue line) and under the optimal policy (dashed red line). It also shows the exogenous state trajectory generated by the linear Markov process in eq. (28) (solid black line). We observe that the trajectory under the learned policy closely matches that of the optimal policy, and both track the mean of the exogenous state.

Section 6.1, we use the knowledge of dynamics in eq. (31) and the true distributions in eq. (28) to compute the true parameters, θ_t^* for $t \in \{1, \dots, T\}$, via the results in Appendix C.1. These are then used to compute the true optimal policy π_t^* as in Corollary 3.2. Finally, we use the true parameters, we compute the regret in eq. (17).

We present our numerical results in Figure 3. The results show that Algorithm 2 is able to learn an effective control policy even when the system matrices A and B are unknown. In particular, the regret $\mathcal{R}(L)$ shown in Figure 3(a) grows sublinearly with the number of episodes L . The average regret shown in Figure 3(b) decreases as L increases. Figure 3(c) presents the norm of the estimation error between the learned and the true parameters, averaged over the episode horizon T , as a function of L . This is expressed as $\|\theta_{1:T}^L - \theta_{1:T}^*\|/T$, where $\theta_{1:T}^L \in \mathbb{R}^{n_\theta \times T}$ is a matrix whose columns correspond to the learned parameters θ_t^L at each time step t with $n_\theta = ((n+m+1)d + (n+m)(n+m+1))/2 = 14$, and similarly for $\theta_{1:T}^*$. We observe that the estimation error decreases with L , indicating that the learned parameters gradually converge to their true values, and consequently that the learned policy approaches the optimal one. Finally, we apply both the learned policy at episode $L = 1500$ and the optimal policy to the system in eq. (31), and compare their corresponding closed-loop state trajectories, as shown in Figure 3(d), alongside the trajectory of the exogenous state s . We observe that the trajectory under the learned policy closely matches that of the optimal policy, showing that Algorithm 2 successfully recovers near-optimal closed-loop behavior despite the unknown system dynamics.

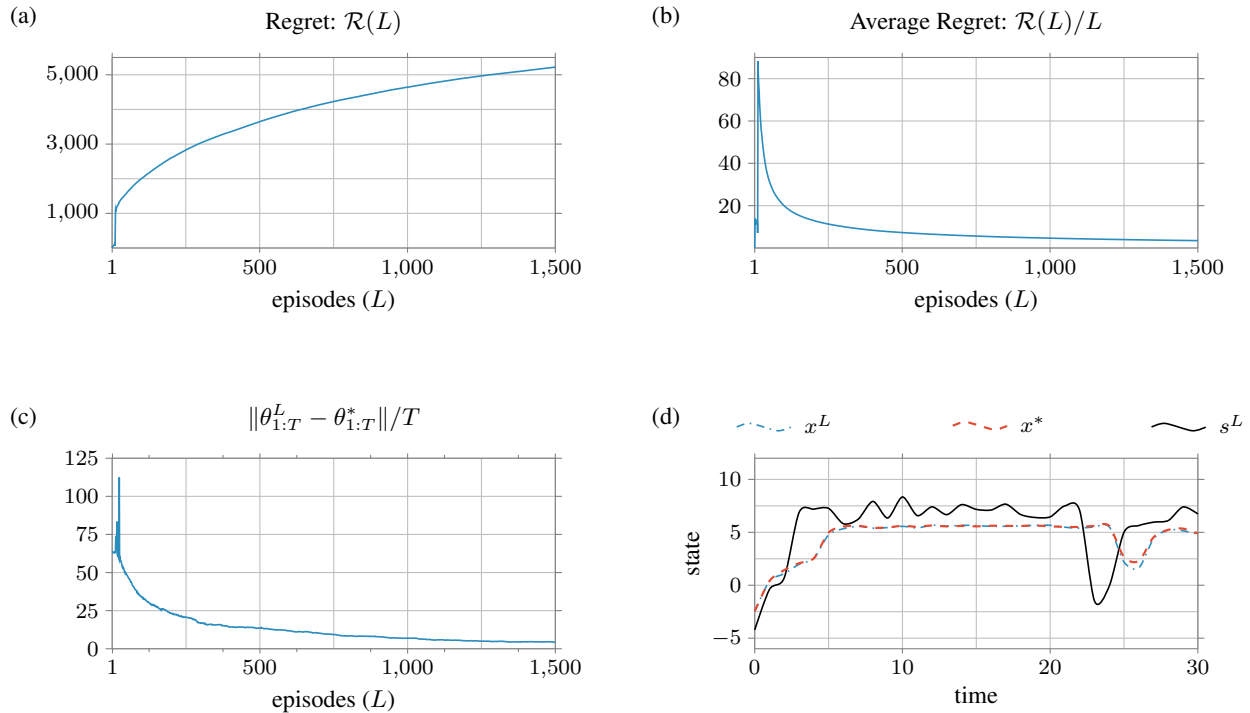


Figure 3: This figure shows the numerical results for the setting in Section 6.3. Panel (a) shows the regret as a function of L , scaling sublinearly with L . Panel (b) shows the average regret as a function of L , and we observe that it converges as L increases. Panel (c) shows the norm of the estimation error between the learned and the true parameters, averaged over the episode horizon T , as a function of L . We observe that the estimation error decreases with L , indicating that the learned policy converges to the optimal one. Panel (d) shows the state trajectory generated by the system in eq. (31) under the learned policy at episode $L = 1500$ (dot-dashed blue line) and under the optimal policy (dashed red line). It also shows the exogenous state trajectory generated by the linear Markov process in eq. (28) (solid black line). We observe that the trajectory under the learned policy closely matches that of the optimal policy, and both track the mean of the exogenous state.

We have also observed that the regret may not converge if the open loop system is unstable (i.e., A has eigenvalues with magnitude greater than 1). Additional enhancements to Algorithm 2 are needed to ensure closed-loop stability. We leave it for future exploration.

7 Conclusion

In this work, we proposed a reinforcement learning framework that unifies linear control systems and feature-based linear Markov models, capturing both deterministic system dynamics and stochastic environmental effects. By leveraging this structure, we derived closed-form expressions for the optimal value function and policy, and introduced a least-squares value iteration algorithm that learns the optimal control policy without requiring explicit model identification or exploration. We provided theoretical guarantees on stability and convergence, and demonstrated the effectiveness of our approach through numerical simulations. We also extended the framework to systems with process noise and to the case of unknown system dynamics. Future directions include establishing stability and convergence guarantees for these extended settings, as well as developing extensions beyond linear system structure.

Broader Impact Statement

This work develops a reinforcement learning framework for controlling linear dynamical systems that interact with stochastic environments, with potential applications in areas such as power systems, robotics, and other cyber-physical systems. The proposed framework may have a positive impact by enabling more adaptive and computationally efficient control under uncertainty, while also incorporating stability considerations that are important in safety-critical applications.

At the same time, like many learning-based control methods, this work could have negative consequences if deployed without sufficient validation and safety safeguards. In particular, the proposed approach is developed under specific modeling assumptions, and its theoretical guarantees are established for the nominal setting. Although we also consider extensions involving Gaussian process noise and unknown system dynamics, these settings are illustrated numerically and are not yet supported by corresponding stability and convergence guarantees. As a result, applying our approach directly to real-world safety-critical systems without additional verification, robustness analysis, and domain-specific safeguards may lead to unsafe or unreliable behavior.

We therefore view this work as a step toward more principled learning-based control, while emphasizing that practical deployment should be accompanied by careful system identification, validation, monitoring, and human oversight.

References

- W. Aangenent, D. Kostic, B. de Jager, R. van de Molengraft, and M. Steinbuch. Data-based optimal control. pp. 1460–1465, Portland, OR, USA, June 2005.
- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, 2011.
- B. D. Anderson and J. B. Moore. *Optimal control: linear quadratic methods*. Courier Corporation, 2007. ISBN 978-0-13-638560-8.
- S. J. Bradtke and A. G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1):33–57, 1996.
- S. J. Bradtke, B. E. Ydstie, and A. G. Barto. Adaptive linear quadratic control using policy iteration. volume 3, pp. 3475–3479. IEEE, 1994.
- J. Bu, A. Mesbahi, M. Fazel, and M. Mesbahi. Lqr through the lens of first order methods: Discrete-time case. *arXiv preprint arXiv:1907.08921*, 2019.
- F. Celi, G. Baggio, and F. Pasqualetti. Closed-form estimates of the LQR gain from finite data. pp. 4016–4021, Cancún, Mexico, December 2022.
- F. Celi, G. Baggio, and F. Pasqualetti. Closed-form and robust formulas for data-driven LQ control. 56, 2023. doi: 10.1016/j.arcontrol.2023.100916.
- G. R. G. da Silva, A. S. Bazanella, C. Lorenzini, and L. Campestri. Data-driven LQR control design. 3(1): 180–185, 2018.
- C. De Persis and P. Tesi. Formulas for data-driven control: Stabilization, optimality and robustness. 65(3): 909–924, 2020.
- S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- F. Dörfler, P. Tesi, and C. De Persis. On the role of regularization in direct data-driven LQR control. pp. 1091–1098, Cancún, Mexico, December 2022. IEEE.
- S. Du, S. Kakade, J. Lee, S. Lovett, G. Mahajan, W. Sun, and R. Wang. Bilinear classes: A structural framework for provable generalization in rl. In *International Conference on Machine Learning*, pp. 2826–2836. PMLR, 2021.

- I. Fatkhullin and B. Polyak. Optimizing static linear feedback: Gradient method. *SIAM Journal on Control and Optimization*, 59(5):3887–3911, 2021.
- M. Fazel, R. Ge, S. Kakade, and M. Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *International Conference on Machine Learning*, pp. 1467–1476, Stockholm, Sweden, 2018.
- F. S. Francisco and M. I. Ribeiro. Q-learning with linear function approximation. In *International Conference on Computational Learning Theory*, pp. 308–322. Springer, 2007.
- N. Jiang, A. Krishnamurthy, A. Agarwal, J. Langford, and R. E. Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1704–1713, 2017.
- C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In Jacob Abernethy and Shivani Agarwal (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 2137–2143. PMLR, 09–12 Jul 2020.
- J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013. doi: 10.1177/0278364913495721.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- H. Mohammadi, A. Zare, M. Soltanolkotabi, and M. R. Jovanović. Global exponential convergence of gradient methods over the nonconvex landscape of the linear quadratic regulator. pp. 7474–7479, Nice, France, Dec. 2019.
- R. Munos. Error bounds for approximate value iteration. In *Proceedings of the National Conference on Artificial Intelligence*, volume 20, pp. 1006, 2005.
- K. Ogata. *Modern Control Engineering*. Instrumentation and controls series. Prentice Hall, 2010. ISBN 9780136156734.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, and et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- W. Sun, N. Jiang, A. Krishnamurthy, A. Agarwal, and J. Langford. Model-based RL in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on learning theory*, pp. 2898–2933. PMLR, 2019.
- R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012. doi: 389–434.
- L. Yang and M. Wang. Sample-optimal parametric q-learning using linearly additive features. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6995–7004, 2019.
- L. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pp. 10746–10756. PMLR, 2020.
- A. Zanette, A. Lazaric, M. Kochenderfer, and E. Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pp. 10978–10989. PMLR, 2020.

A Proof of Theorem 3.1

We first present the following result in which we provide an expression for the optimal greedy policy, π_t^* , and the optimal value function, $V_t^*(x, s)$ under the greedy policy.

Theorem A.1. (optimal policy and value function) Consider the dynamics in eq. (1) and the Markov Process in eq. (2). Let Assumption 2.1 be satisfied. Then, for any $x \in \mathcal{X}$, $s \in \mathcal{S}$, $s' \sim \mathbb{P}(\cdot|s)$, and $t \geq 0$

$$u^*(x, s, t) = \underbrace{K_{x,t}x + K_{h,t}\mathbb{E}[h_{t+1}(s')|s] + K_{s,t}s}_{\pi_t^*(x,s)}, \quad (32)$$

where

$$\begin{aligned} K_{x,t} &= -(R + B^\top G_{t+1}B)^{-1} (B^\top G_{t+1}A + D^\top), \\ K_{h,t} &= -(R + B^\top G_{t+1}B)^{-1} B^\top, \\ K_{s,t} &= -(R + B^\top G_{t+1}B)^{-1} H^\top. \end{aligned} \quad (33)$$

Further,

$$V_t^*(x, s) = x^\top G_t x + 2h_t^\top(s)x + q_t(s), \quad (34)$$

where $G_t \in \mathbb{R}^{n \times n} \succ 0$, $h_t(\cdot) \in \mathbb{R}^n$, and $q_t(\cdot) \in \mathbb{R}$ satisfy

$$G_t = A^\top G_{t+1}A + W - (A^\top G_{t+1}B + D)(R + B^\top G_{t+1}B)^{-1}(B^\top G_{t+1}A + D^\top), \quad (35)$$

$$h_t(s_t) = (A^\top + K_{x,t}^\top B^\top) \mathbb{E}[h_{t+1}(s_{t+1})|s_t] + (F + K_{x,t}^\top H^\top) s_t, \quad (36)$$

$$\begin{aligned} q_t(s_t) &= \mathbb{E}[q_{t+1}(s_{t+1})|s_t] + s_t^\top (M + HK_{s,t})s_t + \mathbb{E}[h_{t+1}^\top(s_{t+1})|s_t] BK_{h,t} \mathbb{E}[h_{t+1}(s_{t+1})|s_t] \\ &\quad + 2s_t^\top HK_{h,t} \mathbb{E}[h_{t+1}(s_{t+1})|s_t], \end{aligned} \quad (37)$$

with $G_T = W$, and $h_T(s_T) = F s_T$ and $q_T(s_T) = s_T^\top M s_T$.

Proof. We prove our claim by induction. For notational convenience, we drop the time index from the states and inputs inside the expressions and arguments of $c(\cdot)$, $V_t^*(\cdot)$, and $Q_t^*(\cdot)$ for $t \geq 0$, where we use x , s , u , x' , and s' to denote x_t , s_t , u_t , x_{t+1} , and s_{t+1} , respectively. At $t = T - 1$,

$$\begin{aligned} Q_{T-1}^*(x, s, u) &= c(x, s, u) + \mathbb{E}_{s' \sim \mathbb{P}_t(s'|s)} [V_T^*(x', s')|x, s, u] \\ &= x^\top (W + A^\top W A)x + u^\top (R + B^\top W B)u + 2(s^\top F^\top + \mathbb{E}[s'|s]^\top F^\top A)x \\ &\quad + 2 \left(x^\top (D + A^\top W B) + s^\top H + \mathbb{E}[s'|s]^\top F^\top B \right) u + s^\top M s + \mathbb{E}[s'^\top M s'|s], \end{aligned} \quad (38)$$

where we used the fact that $V_T^*(x_T, s_T) = c(x_T, s_T, u_T)$ and $u_T = 0$. Taking the derivative of Q_{T-1}^* with respect to u

$$\frac{\partial Q_{T-1}^*(x, s, u)}{\partial u} = 2(R + B^\top W B)u + 2(B^\top W A + D^\top)x + 2(B^\top F \mathbb{E}[s'|s] + H^\top s).$$

Setting the above derivative to zero and solving for u , we get

$$u_{T-1}^* = -(R + B^\top W B)^{-1}(B^\top W A + D^\top)x_{T-1} - (R + B^\top W B)^{-1}(B^\top F \mathbb{E}[s_T|s_{T-1}] + H^\top s_{T-1}), \quad (39)$$

which is the minimizer of eq. (38). We substitute eq. (39) in eq. (38),

$$V_{T-1}^*(x, s) = x^\top G_{T-1}x + 2h_{T-1}^\top(s)x + q_{T-1}(s), \quad (40)$$

where G_{T-1} , $h_{T-1}(s)$, and $q_{T-1}(s)$ are as in Theorem A.1 for $t = T - 1$. Suppose for $t = k + 1$,

$$V_{k+1}^*(x, s) = x^\top G_{k+1}x + 2h_{k+1}^\top(s)x + q_{k+1}(s),$$

where G_{k+1} , $h_{k+1}(s)$, and $q_{k+1}(s)$ are as in Theorem A.1 for $t = k + 1$. Then we have,

$$\begin{aligned} Q_k^*(x, s, u) &= c(x, s, u) + \mathbb{E}_{s' \sim \mathbb{P}_t(s'|s)} [V_{k+1}^*(x', s') | x, s, u] \\ &= x^\top (W + A^\top G_{k+1} A) x + u^\top (R + B^\top G_{k+1} B) u + 2(s^\top F^\top + \mathbb{E}[h_{k+1}(s') | s]^\top A) x \\ &\quad + 2 \left(x^\top (D + A^\top G_{k+1} B) + s^\top H + \mathbb{E}[h_{k+1}(s') | s]^\top B \right) u + s^\top M s + \mathbb{E}[q_{k+1}(s') | s], \end{aligned} \quad (41)$$

Taking the derivative of Q_k^* with respect to u

$$\frac{\partial Q_k^*(x, s, u)}{\partial u} = 2(R + B^\top G_{k+1} B) u + 2(B^\top G_{k+1} A + D^\top) x + 2(B^\top \mathbb{E}[h_{k+1}(s') | s] + H^\top s).$$

Setting the above derivative to zero and solving for u , we get

$$\begin{aligned} u_k^* &= - (R + B^\top G_{k+1} B)^{-1} (B^\top G_{k+1} A + D^\top) x_k - (R + B^\top G_{k+1} B)^{-1} H^\top s_k \\ &\quad - (R + B^\top G_{k+1} B)^{-1} B^\top \mathbb{E}[h_{k+1}(s(k+1)) | s_k], \end{aligned} \quad (42)$$

which is the minimizer of eq. (41). We substitute eq. (42) in eq. (41),

$$V_k^*(x, s) = x^\top G_k x + 2h_k(s)^\top x + q_k(s),$$

where G_k , $h_k(s)$, and $q_k(s)$ are as in Theorem A.1 for $t = k$. This completes the proof. \square

Proof of Theorem 3.1: For notational convenience, we drop the time index from the states and inputs inside the expressions and arguments of $c(\cdot)$, $V_t^*(\cdot)$, and $Q_t^*(\cdot)$ for $t \geq 0$, where we use x , s , u , x' , and s' to denote x_t , s_t , u_t , x_{t+1} , and s_{t+1} , respectively. Using Theorem A.1, we write

$$\begin{aligned} Q_t^*(x, s, u) &= c(x, s, u) + \mathbb{E}_{s' \sim \mathbb{P}_t(s'|s)} [V_{t+1}^*(x', s') | x, s, u] \\ &= c(x, s, u) + \int_{\mathcal{S}} V_{t+1}^*(Ax + Bu, s') \mathbb{P}_t(ds' | s) \\ &\stackrel{(a)}{=} c(x, s, u) + \int_{\mathcal{S}} V_{t+1}^*(Ax + Bu, s') \phi(s)^\top \mu_t(ds') \\ &\stackrel{(b)}{=} c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \int_{\mathcal{S}} (2h_{t+1}(s')^\top (Ax + Bu) + q_{t+1}(s')) \sum_{i=1}^d \phi_i(s) \mu_{i,t}(ds') \\ &= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \int_{\mathcal{S}} q_{t+1}(s') \mu_{i,t}(ds') \\ &\quad + 2 \sum_{i=1}^d \left(\phi_i(s) \int_{\mathcal{S}} h_{t+1}(s')^\top \mu_{i,t}(ds') \right) (Ax + Bu) \\ &= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \underbrace{\mathbb{E}_{\mu_{i,t}} [q_{t+1}(s')]}_{\bar{q}_{i,t+1}} \\ &\quad + 2 \sum_{i=1}^d \left(\phi_i(s) \underbrace{\mathbb{E}_{\mu_{i,t}} [h_{t+1}(s')^\top]}_{\bar{h}_{i,t+1}^\top} \right) (Ax + Bu), \end{aligned}$$

where in step (a) we have used Assumption 2.1, and in step (b) we have used Theorem A.1. \blacksquare

B Least-squares value iteration

We formulate the regularized least squares regression and derive its solution presented in lines 6-7 of Algorithm 1. We begin by using the notation in eq. (9) to derive the expression of the parametrized state-action value

function, Q_t and the corresponding parametrized optimal greedy policy in eq. (10). Using the expression of Q_t in Theorem 3.1, we can write

$$\begin{aligned}
Q_t(x, s, u) &= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \left(2 (Ax + Bu)^\top \bar{h}_{i,t+1} + \bar{q}_{i,t+1} \right) \\
&= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) \underbrace{\left[2 (Ax + Bu)^\top \quad 1 \right]}_{y(x,u)^\top} \underbrace{\begin{bmatrix} \bar{h}_{i,t+1} \\ \bar{q}_{i,t+1} \end{bmatrix}}_{\theta_{i,t+1}} \\
&= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \sum_{i=1}^d \phi_i(s) y(x, u)^\top \theta_{i,t+1} \\
&= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \underbrace{\phi(s)^\top \left(I_d \otimes y(x, u)^\top \right)}_{Y(x,u)} \underbrace{\begin{bmatrix} \theta_{1,t+1} \\ \vdots \\ \theta_{d,t+1} \end{bmatrix}}_{\theta_{t+1}} \\
&= c(x, s, u) + (Ax + Bu)^\top G_{t+1} (Ax + Bu) + \phi(s)^\top Y(x, u) \theta_{t+1}.
\end{aligned} \tag{43}$$

Next, using the notation in eq. (9), we can write

$$\begin{aligned}
\bar{h}_{i,t+1} &= \underbrace{\begin{bmatrix} I_n & 0_{n \times 1} \end{bmatrix}}_Z \underbrace{\begin{bmatrix} \bar{h}_{i,t+1} \\ \bar{q}_{i,t+1} \end{bmatrix}}_{\theta_{i,t+1}} = Z \theta_{i,t+1}, \quad i \in \{1, \dots, d\}, \\
\bar{h}_{t+1} &= \begin{bmatrix} Z & 0 & \cdots & 0 \\ 0 & Z & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & Z \end{bmatrix} \begin{bmatrix} \theta_{1,t+1} \\ \vdots \\ \theta_{d,t+1} \end{bmatrix} = (I_d \otimes Z) \theta_{t+1}.
\end{aligned} \tag{44}$$

Similarly, we can write

$$\begin{aligned}
\bar{q}_{i,t+1} &= \underbrace{\begin{bmatrix} 0_{1 \times n} & 1 \end{bmatrix}}_{\bar{Z}} \underbrace{\begin{bmatrix} \bar{h}_{i,t+1} \\ \bar{q}_{i,t+1} \end{bmatrix}}_{\theta_{i,t+1}} = \bar{Z} \theta_{i,t+1}, \quad i \in \{1, \dots, d\}, \\
\bar{q}_{t+1} &= \begin{bmatrix} \bar{Z} & 0 & \cdots & 0 \\ 0 & \bar{Z} & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \bar{Z} \end{bmatrix} \begin{bmatrix} \theta_{1,t+1} \\ \vdots \\ \theta_{d,t+1} \end{bmatrix} = (I_d \otimes \bar{Z}) \theta_{t+1}.
\end{aligned} \tag{45}$$

Next, from the expression of u_t in Corollary 3.2, we write

$$\begin{aligned}
u_t(x, s) &= K_{x,t} x + K_{s,t} s + K_{h,t} \sum_{i=1}^d \phi_i(s) \bar{h}_{i,t+1} \\
&= K_{x,t} x + K_{s,t} s + K_{h,t} \left[\phi_1(s) I_n \quad \cdots \quad \phi_d(s) I_n \right] \begin{bmatrix} \bar{h}_{1,t+1} \\ \vdots \\ \bar{h}_{d,t+1} \end{bmatrix} \\
&= K_{x,t} x + K_{s,t} s + K_{h,t} \left(\phi(s)^\top \otimes I_n \right) \bar{h}_{t+1} \\
&\stackrel{(a)}{=} K_{x,t} x + K_{s,t} s + \left(\phi(s)^\top \otimes I_n \right) (I_d \otimes Z) \theta_{t+1} \\
&= K_{x,t} x + K_{s,t} s + \left(\phi(s)^\top \otimes Z \right) \theta_{t+1},
\end{aligned} \tag{46}$$

where in step (a) we have used eq. (44). We define the Bellman target at time t as

$$g_t(x, s, u) = c(x, s, u) + \min_v \widehat{Q}_{t+1}(x', s', v), \quad (47)$$

where x' and s' denote the states resulting from taking action u in states x and s , and $\widehat{Q}_{t+1}(x', s', v)$ is the estimate of the state-action value function at time $t+1$. We re-write eq. (48) as

$$\begin{aligned} g_t(x, s, u) &= c(x, s, u) + \widehat{V}_{t+1}^*(x', s') \\ &\stackrel{(b)}{=} c(x, s, u) + x'^T G_{t+1} x' + 2\widehat{h}_{t+1}^T(s') x' + \widehat{q}_{t+1}(s'), \end{aligned} \quad (48)$$

where in step (b), we have used Theorem A.1. For notational convenience, let $X_1(t) = A^T + K_{x,t}^T B^T$, $X_2(t) = F + K_{x,t}^T H^T$, $Y_1(t) = M + H K_{s,t}$, $Y_2(t) = B K_{h,t}$, and $Y_3(t) = H K_{h,t}$. Using eq. (36) and eq. (37), we re-write \widehat{h}_{t+1} and \widehat{q}_{t+1} in eq. (48) as

$$\begin{aligned} \widehat{h}_{t+1}(s_{t+1}) &= X_1(t+1) \mathbb{E}[h_{t+2}(s_{t+2}) | s_{t+1}] + X_2(t+1) s_{t+1}, \\ &\stackrel{(c)}{=} X_1(t+1) \left(\phi(s_{t+1})^T \otimes Z \right) \widehat{\theta}_{t+2} + X_2(t+1) s_{t+1}, \end{aligned} \quad (49)$$

$$\begin{aligned} \widehat{q}_{t+1}(s_{t+1}) &= \mathbb{E}[q_{t+2}(s_{t+2}) | s_{t+1}] + s_{t+1}^T Y_1(t+1) s_{t+1} + \mathbb{E}[h_{t+2}^T(s_{t+2}) | s_{t+1}] Y_2(t+1) \mathbb{E}[h_{t+2}(s_{t+2}) | s_{t+1}] \\ &\quad + 2s_{t+1}^T Y_3(t+1) \mathbb{E}[h_{t+2}(s_{t+2}) | s_{t+1}], \\ &\stackrel{(d)}{=} \phi(s_{t+1})^T (I_d \otimes \overline{Z}) \widehat{\theta}_{t+2} + s_{t+1}^T Y_1(t+1) s_{t+1} + \widehat{\theta}_{t+2}^T (\phi(s_{t+1}) \otimes Z^T) Y_2(t+1) \left(\phi(s_{t+1})^T \otimes Z \right) \widehat{\theta}_{t+2} \\ &\quad + 2s_{t+1}^T Y_3(t+1) \left(\phi(s_{t+1})^T \otimes Z \right) \widehat{\theta}_{t+2}, \end{aligned} \quad (50)$$

where in steps (c) and (d) we have used eq. (44) and eq. (45), respectively. The temporal difference (TD) error is written as

$$\begin{aligned} \varepsilon_t(x, s, u) &= g_t(x, s, u) - \widehat{Q}_t(x, s, u) \\ &= c(x, s, u) + \min_v \widehat{Q}_{t+1}(x', s', v) - \widehat{Q}_t(x, s, u) \\ &\stackrel{(d)}{=} c(x, s, u) + x'^T G_{t+1} x' + 2\widehat{h}_{t+1}^T(s') x' + \widehat{q}_{t+1}(s') - c(x, s, u) - (Ax + Bu)^T G_{t+1} (Ax + Bu) \\ &\quad - \phi(s)^T Y(x, u) \widehat{\theta}_{t+1} \\ &= 2\widehat{h}_{t+1}^T(s') x' + \widehat{q}_{t+1}(s') - \phi(s)^T Y(x, u) \widehat{\theta}_{t+1}, \end{aligned} \quad (51)$$

where in step (d) we have used eq. (43) and eq. (48). The TD error $\varepsilon_t(x, s, u)$ in eq. (51) captures the discrepancy between the Bellman target and the current estimate of the Q -function. In Least-Squares Value Iteration (LSVI) in Algorithm 1, we minimize the squared TD error over the dataset, eq. (11), collected up to episode $\ell-1$, to obtain an updated estimate of the Q -function. Specifically, the parameters $\widehat{\theta}_t$ at episode ℓ denoted by θ_t^ℓ is obtained by solving the following regularized least-squares problem

$$\begin{aligned} \theta_{t+1}^\ell &= \arg \min_{\theta} \underbrace{\sum_{j=1}^{\ell-1} \varepsilon_t(x^j, s^j, u^j)^2}_J + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \sum_{j=1}^{\ell-1} \left(\phi(s^j)^T Y(x^j, u^j) \widehat{\theta}_{t+1} - 2\widehat{h}_{t+1}^T(s^j) x'^j - \widehat{q}_{t+1}(s^j) \right)^2 + \lambda \|\theta\|_2^2. \end{aligned} \quad (52)$$

Taking the derivative of eq. (52) with respect to θ , we get

$$\frac{\partial J}{\partial \theta} = 2 \sum_{j=1}^{\ell-1} Y^T(x^j, u^j) \phi(s^j) \left(\phi(s^j)^T Y(x^j, u^j) \theta - 2\widehat{h}_{t+1}^T(s^j) x'^j - \widehat{q}_{t+1}(s^j) \right) + 2\lambda \theta.$$

Setting the above derivative to zero and solving for θ , we get

$$\begin{aligned}\theta_{t+1}^\ell &= \Lambda_t^{-1} \sum_{j=1}^{\ell-1} Y^\top(x^j, u^j) \phi(s^j) (2\widehat{h}_{t+1}^\top(s^j) x^j + \widehat{q}_{t+1}(s^j)), \\ \Lambda_t &= \sum_{j=1}^{\ell-1} Y(x^j, u^j)^\top \phi(s^j) \phi(s^j)^\top Y(x^j, u^j) + \lambda I_{d(n+1)}.\end{aligned}\tag{53}$$

In Algorithm 1 we computed \widehat{h}_{t+1} and \widehat{q}_{t+1} at episode ℓ as in eq. (49) and eq. (50), respectively. These quantities are obtained using the updated parameter θ_{t+2}^ℓ from the previous iteration of the backward-in-time weight update loop (lines 5-11 in Algorithm 1). In particular, we have $\widehat{h}_{t+1}(\cdot) = h_{t+1}^\ell(\cdot)$ and $\widehat{q}_{t+1}(\cdot) = q_{t+1}^\ell(\cdot)$.

C True weights \bar{h}_t and \bar{q}_t

Throughout this Appendix, we use the following notation,

$$\begin{aligned}X_1(t) &= A^\top + K_{x,t}^\top B^\top, & X_2(t) &= F + K_{x,t}^\top H^\top, \\ Y_1(t) &= M + HK_{s,t}, & Y_2(t) &= BK_{h,t}, & Y_3(t) &= HK_{h,t}, \\ \Phi_t &= \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [\phi(s_t)^\top] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [\phi(s_t)^\top] \end{bmatrix}, & \text{and } \bar{m}_t &= \begin{bmatrix} \bar{m}_{1,t} \\ \vdots \\ \bar{m}_{d,t} \end{bmatrix} = \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [s_t] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [s_t] \end{bmatrix},\end{aligned}\tag{54}$$

for $t \in \{0, \dots, T\}$. In addition, we define for $t \in \{0, \dots, T\}$

$$\|X_1(t)\| \leq \bar{X}_1, \quad \|X_2(t)\| \leq \bar{X}_2, \quad \|Y_1(t)\| \leq \bar{Y}_1, \quad \|Y_2(t)\| \leq \bar{Y}_2, \quad \|Y_3(t)\| \leq \bar{Y}_3.\tag{55}$$

C.1 Closed-form expressions of \bar{h}_t and \bar{q}_t

In this Appendix, we derive closed-form expressions for the true parameters $\bar{h}_t = \mathbb{E}_{\mu_t} [h_t(s_t)]$ and $\bar{q}_t = \mathbb{E}_{\mu_t} [q_t(s_t)]$.

Theorem C.1. (closed-form expressions for the true \bar{h}_t and \bar{q}_t) Consider the dynamics in eq. (1) and the Markov Process in eq. (2). Let Assumption 2.1 and Assumption 2.2 be satisfied. Let $\bar{h}_t = \mathbb{E}_{\mu_t} [h_t(s_t)]$ and $\bar{q}_t = \mathbb{E}_{\mu_t} [q_t(s_t)]$ for $t \in \{0, \dots, T\}$, where $h_t(\cdot)$ and $q_t(\cdot)$ are as in eq. (36) and eq. (37), respectively. Then, for $t \in \{0, \dots, T-1\}$

$$\begin{aligned}\bar{h}_t &= (\Phi_t \otimes X_1(t)) \bar{h}_{t+1} + (I_d \otimes X_2(t)) \bar{m}_t, \\ \bar{q}_t &= \Phi_t \bar{q}_{t+1} + \mathbb{E}_{\mu_t} [s_t^\top Y_1(t) s_t] + \begin{bmatrix} \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{1,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \\ \vdots \\ \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{d,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \end{bmatrix} + 2\mathbb{E}_{\mu_t} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)] \bar{h}_{t+1},\end{aligned}$$

with $\bar{h}_T = F \mathbb{E}_{\mu_T} [s_T]$ and $\bar{q}_T = \mathbb{E}_{\mu_T} [s_T^\top M s_T]$, where $X_1(t)$, $X_2(t)$, $Y_1(t)$, $Y_2(t)$, $Y_3(t)$, Φ_t , and \bar{m}_t are as in eq. (54).

Proof. We re-write equation eq. (36) as

$$h_t(s_t) = X_1(t) \mathbb{E} [h_{t+1}(s_{t+1}) | s_t] + X_2(t) s_t.\tag{56}$$

Taking the expectation of both sides with respect to $\mu_{i,t}$ for each $i \in \{1, \dots, d\}$, we get

$$\begin{aligned}
\mathbb{E}_{\mu_{i,t}} [h_t(s_t)] &= X_1(t) \mathbb{E}_{\mu_{i,t}} [\mathbb{E}[h_{t+1}(s_{t+1}) | s_t]] + X_2(t) \mathbb{E}_{\mu_{i,t}} [s_t] \\
&= X_1(t) \mathbb{E}_{\mu_{i,t}} \left[\sum_{j=1}^d \phi_j(s_t) \mathbb{E}_{\mu_{j,t+1}} [h_{t+1}(s_{t+1})] \right] + X_2(t) \bar{m}_{i,t} \\
&= X_1(t) \sum_{j=1}^d \mathbb{E}_{\mu_{i,t}} [\phi_j(s_t)] \mathbb{E}_{\mu_{j,t+1}} [h_{t+1}(s_{t+1})] + X_2(t) \bar{m}_{i,t} \\
&= X_1(t) (\mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top] \otimes I_n) \bar{h}_{t+1} + X_2(t) \bar{m}_{i,t}.
\end{aligned} \tag{57}$$

By noting that

$$\bar{h}_t = \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [h_t(s_t)] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [h_t(s_t)] \end{bmatrix},$$

and denoting

$$\Phi_t = \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [\phi(s_t)^\top] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [\phi(s_t)^\top] \end{bmatrix}, \quad \text{and} \quad \bar{m}_t = \begin{bmatrix} \bar{m}_{1,t} \\ \vdots \\ \bar{m}_{d,t} \end{bmatrix},$$

we can write

$$\begin{aligned}
\bar{h}_t &= (I_d \otimes X_1(t)) (\Phi_t \otimes I_n) \bar{h}_{t+1} + (I_d \otimes X_2(t)) \bar{m}_t \\
&= (\Phi_t \otimes X_1(t)) \bar{h}_{t+1} + (I_d \otimes X_2(t)) \bar{m}_t.
\end{aligned} \tag{58}$$

Next, we re-write equation eq. (37) as

$$\begin{aligned}
q_t(s_t) &= \mathbb{E}[q_{t+1}(s_{t+1}) | s_t] + s_t^\top Y_1(t) s_t + \mathbb{E}[h_{t+1}^\top(s_{t+1}) | s_t] Y_2(t) \mathbb{E}[h_{t+1}(s_{t+1}) | s_t] \\
&\quad + 2s_t^\top Y_3(t) \mathbb{E}[h_{t+1}(s_{t+1}) | s_t].
\end{aligned} \tag{59}$$

Taking the expectation of both sides with respect to $\mu_{i,t}$ for each $i \in \{1, \dots, d\}$, we get

$$\begin{aligned}
\mathbb{E}_{\mu_{i,t}} [q_t(s_t)] &= \mathbb{E}_{\mu_{i,t}} [\mathbb{E}[q_{t+1}(s_{t+1}) | s_t]] + \mathbb{E}_{\mu_{i,t}} [s_t^\top Y_1(t) s_t] + \mathbb{E}_{\mu_{i,t}} [\mathbb{E}[h_{t+1}^\top(s_{t+1}) | s_t] Y_2(t) \mathbb{E}[h_{t+1}(s_{t+1}) | s_t]] \\
&\quad + 2\mathbb{E}_{\mu_{i,t}} [s_t^\top Y_3(t) \mathbb{E}[h_{t+1}(s_{t+1}) | s_t]].
\end{aligned} \tag{60}$$

We start with,

$$\begin{aligned}
\mathbb{E}_{\mu_{i,t}} [\mathbb{E}[q_{t+1}(s_{t+1}) | s_t]] &= \mathbb{E}_{\mu_{i,t}} \left[\sum_{j=1}^d \phi_j(s_t) \mathbb{E}_{\mu_{j,t+1}} [q_{t+1}(s_{t+1})] \right] \\
&= \mathbb{E}_{\mu_{i,t}} \left[\sum_{j=1}^d \phi_j(s_t) \right] \mathbb{E}_{\mu_{j,t+1}} [q_{t+1}(s_{t+1})] \\
&= \underbrace{[\mathbb{E}_{\mu_{i,t}} [\phi_1(s_t)] \cdots \mathbb{E}_{\mu_{i,t}} [\phi_d(s_t)]]}_{\mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top]} \underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t+1}} [q_{t+1}(s_{t+1})] \\ \vdots \\ \mathbb{E}_{\mu_{d,t+1}} [q_{t+1}(s_{t+1})] \end{bmatrix}}_{\bar{q}_{t+1}}.
\end{aligned} \tag{61}$$

Next we have,

$$\begin{aligned}
& \mathbb{E}_{\mu_{i,t}} \left[\mathbb{E} \left[h_{t+1}^\top (s_{t+1}) | s_t \right] Y_2(t) \mathbb{E} \left[h_{t+1} (s_{t+1}) | s_t \right] \right] \\
&= \mathbb{E}_{\mu_{i,t}} \left[\sum_{j=1}^d \phi_j(s_t) \mathbb{E}_{\mu_{j,t+1}} \left[h_{t+1}(s_{t+1})^\top \right] Y_2(t) \sum_{k=1}^d \phi_k(s_t) \mathbb{E}_{\mu_{k,t+1}} \left[h_{t+1}(s_{t+1}) \right] \right] \\
&= \mathbb{E}_{\mu_{i,t}} \left[\underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t+1}} [h_{t+1}(s_{t+1})] \\ \vdots \\ \mathbb{E}_{\mu_{d,t+1}} [h_{t+1}(s_{t+1})] \end{bmatrix}}_{\bar{h}_{t+1}}^\top \underbrace{\begin{bmatrix} \phi_1(s_t) I_n \\ \vdots \\ \phi_d(s_t) I_n \end{bmatrix}}_{\phi(s_t) \otimes I_n} Y_2(t) \underbrace{\begin{bmatrix} \phi_1(s_t) I_n \\ \vdots \\ \phi_d(s_t) I_n \end{bmatrix}}_{\phi(s_t) \otimes I_n}^\top \underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t+1}} [h_{t+1}(s_{t+1})] \\ \vdots \\ \mathbb{E}_{\mu_{d,t+1}} [h_{t+1}(s_{t+1})] \end{bmatrix}}_{\bar{h}_{t+1}} \right] \quad (62) \\
&= \bar{h}_{t+1}^\top \mathbb{E}_{\mu_{i,t}} \left[(\phi(s_t) \otimes I_n) Y_2(t) (\phi(s_t)^\top \otimes I_n) \right] \bar{h}_{t+1} \\
&= \bar{h}_{t+1}^\top \mathbb{E}_{\mu_i} \left[(\phi(s_t) \otimes I_n) (1 \otimes Y_2(t)) (\phi(s_t)^\top \otimes I_n) \right] \bar{h}_{t+1} \\
&= \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_i} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1}.
\end{aligned}$$

Finally, we have

$$\begin{aligned}
\mathbb{E}_{\mu_{i,t}} \left[s_t^\top Y_3(t) \mathbb{E} \left[h_{t+1} (s_{t+1}) | s_t \right] \right] &= \mathbb{E}_{\mu_{i,t}} \left[s_t^\top Y_3(t) \sum_{j=1}^d \phi_j(s_t) \mathbb{E}_{\mu_{j,t+1}} \left[h_{t+1}(s_{t+1}) \right] \right] \\
&= \mathbb{E}_{\mu_{i,t}} \left[s_t^\top Y_3(t) (\phi(s_t)^\top \otimes I_n) \right] \bar{h}_{t+1} \quad (63) \\
&= \mathbb{E}_{\mu_{i,t}} \left[(1 \otimes s_t^\top Y_3(t)) (\phi(s_t)^\top \otimes I_n) \right] \bar{h}_{t+1} \\
&= \mathbb{E}_{\mu_{i,t}} \left[(\phi(s_t)^\top \otimes s_t^\top Y_3(t)) \right] \bar{h}_{t+1}.
\end{aligned}$$

Substituting eq. (61), eq. (62), and eq. (63) in eq. (60), we get

$$\begin{aligned}
\mathbb{E}_{\mu_{i,t}} [q_t(s_t)] &= \mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top] \bar{q}_{t+1} + \mathbb{E}_{\mu_{i,t}} [s_t^\top Y_1(t) s_t] \\
&\quad + \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{i,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \quad (64) \\
&\quad + 2 \mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)] \bar{h}_{t+1}.
\end{aligned}$$

Then, we can write

$$\begin{aligned}
\underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [q_t(s_t)] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [q_t(s_t)] \end{bmatrix}}_{\bar{q}_t} &= \underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [\phi(s_t)^\top] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [\phi(s_t)^\top] \end{bmatrix}}_{\Phi_t} \underbrace{\begin{bmatrix} \mathbb{E}_{\mu_{1,t+1}} [q_{t+1}(s_{t+1})] \\ \vdots \\ \mathbb{E}_{\mu_{d,t+1}} [q_{t+1}(s_{t+1})] \end{bmatrix}}_{\bar{q}_{t+1}} + \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [s_t^\top Y_1(t) s_t] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [s_t^\top Y_1(t) s_t] \end{bmatrix} \\
&\quad + \begin{bmatrix} \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{1,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \\ \vdots \\ \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{d,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \end{bmatrix} + 2 \begin{bmatrix} \mathbb{E}_{\mu_{1,t}} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)] \\ \vdots \\ \mathbb{E}_{\mu_{d,t}} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)] \end{bmatrix} \bar{h}_{t+1}. \quad (65)
\end{aligned}$$

Finally, from Theorem A.1, we have $\bar{h}_T = \mathbb{E}_{\mu_T} [h_T(s_T)] = F \mathbb{E}_{\mu_T} [s_T]$, and $\bar{q}_t = \mathbb{E}_{\mu_T} [q_T(s_T)] = \mathbb{E}_{\mu_T} [s_T^\top M s_T]$. \square

C.2 Upper bounds on $\|\bar{h}_t\|$ and $\|\bar{q}_t\|$

Theorem C.2. (upper bounds on the true \bar{h}_t and \bar{q}_t) Consider the dynamics in eq. (1) and the Markov Process in eq. (2). Let Assumption 2.1 and Assumption 2.2 be satisfied. Let \bar{h}_t and \bar{q}_t be as in Theorem C.1 for $t \in \{0, \dots, T\}$. Then,

$$\begin{aligned}
\|\bar{h}_t\| &\leq \|F\| \delta_s \alpha \rho^{T-t} + \frac{\bar{X}_2 \delta_s \alpha \sqrt{d}}{1 - \rho}, \\
\|\bar{q}_t\| &\leq \delta_s^2 \|M\| + \delta_s^2 \bar{Y}_1 \sqrt{d} + \frac{\bar{Y}_2 \|\bar{h}_{t+1}\|^2}{\sqrt{d}} + \delta_s \bar{Y}_3 \|\bar{h}_{t+1}\|,
\end{aligned}$$

for $t \in \{0, \dots, T\}$ with $\|\bar{h}_{T+1}\| = 0$, where $\alpha > 0$, $0 < \rho < 1$ are constants, and $\bar{X}_1, \bar{X}_2, \bar{Y}_1, \bar{Y}_2$, and \bar{Y}_3 are as in eq. (55).

Proof. From eq. (58), we have

$$\bar{h}_t = (\Phi_t \otimes X_1(t))\bar{h}_{t+1} + (I_d \otimes X_2(t))\bar{m}_t. \quad (66)$$

Let $\Xi(t_1, t_2) = \prod_{i=t_2-1}^{t_1} \Phi_i \otimes X_1(i)$ with $t_2 > t_1$. Then, given \bar{h}_T , we can write

$$\bar{h}_t = \Xi(t, T)\bar{h}_T + \sum_{j=T-1}^t \Xi(t, j) (I_d \otimes X_2(j))\bar{m}_j. \quad (67)$$

Then, we can write

$$\|\bar{h}_t\| \leq \|\Xi(t, T)\| \|\bar{h}_T\| + \sum_{j=T-1}^t \|\Xi(t, j)\| \|(I_d \otimes X_2(j))\| \|\bar{m}_j\|. \quad (68)$$

Now bound each term separately. For $t_2 > t_1$, we have

$$\Xi(t_1, t_2) = \prod_{i=t_2-1}^{t_1} \Phi_i \otimes X_1(i) = \left(\prod_{i=t_2-1}^{t_1} \Phi_i \right) \otimes \left(\prod_{i=t_2-1}^{t_1} X_1(i) \right). \quad (69)$$

Notice that, from eq. (36) and eq. (26), we have $X_1(i) = A^\top + K_{i,x}^\top B^\top = A_c(i)^\top$. Then, we can write

$$\prod_{i=t_2-1}^{t_1} X_1(i) = \prod_{i=t_2-1}^{t_1} A_c(i)^\top = \left(\prod_{i=t_1}^{t_2-1} A_c(i) \right)^\top. \quad (70)$$

Then, noting that $\|\cdot^\top\| = \|\cdot\|$, we can upper bound eq. (69) as

$$\|\Xi(t_1, t_2)\| \leq \left(\prod_{i=t_2-1}^{t_1} \|\Phi_i\| \right) \left(\left\| \prod_{i=t_1}^{t_2-1} A_c(i) \right\| \right). \quad (71)$$

From (Celi et al., 2022, Lemma B.1), we have $\left\| \prod_{i=t_1}^{t_2-1} A_c(i) \right\| \leq \alpha \rho^{t_2-t_1}$, where $\alpha > 0$ and $\rho \in (0, 1)$ are constants. Further, using Assumption 2.1, we have for any $i \leq 0$

$$\begin{aligned} \|\Phi_i\|_2 \leq \|\Phi_i\|_F &= \sqrt{\text{tr}\left((\mathbb{E}_{\mu_i}[\phi^\top(s_i)]) (\mathbb{E}_{\mu_i}[\phi^\top(s_i)])^\top \right)} \\ &= \sqrt{\sum_{j=1}^d (\mathbb{E}_{\mu_{j,i}}[\phi^\top(s_i)]) (\mathbb{E}_{\mu_{j,i}}[\phi^\top(s_i)])^\top} \\ &= \sqrt{\sum_{j=1}^d \|\mathbb{E}_{\mu_{j,i}}[\phi^\top(s_i)]\|_2^2} \\ &\leq \sqrt{\sum_{j=1}^d \frac{1}{d}} = 1. \end{aligned} \quad (72)$$

Hence, we can re-write eq. (71) as

$$\|\Xi(t_1, t_2)\| \leq \alpha \rho^{t_2-t_1}. \quad (73)$$

From Theorem A.1, we have $\bar{h}_T = \mathbb{E}_{\mu_T} [h_T(s_T)] = F\mathbb{E}_{\mu_T} [s_T]$. Then, we have

$$\|\bar{h}_T\| \leq \|F\| \|\mathbb{E}_{\mu_T} [s_T]\| \stackrel{(a)}{\leq} \|F\| \mathbb{E}_{\mu_T} [\|s_T\|] \stackrel{(b)}{\leq} \|F\| \delta_s, \quad (74)$$

where in step (a) we have used the Jensen's inequality, and in step (b) we have used Assumption 2.1. Next, we have

$$\|\bar{m}_t\| = \|\mathbb{E}_{\mu_t} [s_t]\| = \sqrt{\sum_{i=1}^d (\mathbb{E}_{\mu_{i,t}} [s_t])^\top (\mathbb{E}_{\mu_{i,t}} [s_t])} = \sqrt{\sum_{i=1}^d \|\mathbb{E}_{\mu_{i,t}} [s_t]\|^2} \leq \sqrt{d\delta_s^2} = \delta_s \sqrt{d}. \quad (75)$$

Let $\|X_2(t)\| \leq \bar{X}_2$ for $t \in \{0, \dots, T-1\}$. Then, using eq. (73), eq. (74), and eq. (75) we can write eq. (68) as

$$\begin{aligned} \|\bar{h}_t\| &\leq \|F\| \delta_s \alpha \rho^{T-t} + \bar{X}_2 \delta_s \alpha \sqrt{d} \sum_{j=T-1}^t \rho^{j-t} \\ &\stackrel{(c)}{=} \|F\| \delta_s \alpha \rho^{T-t} + \bar{X}_2 \delta_s \alpha \sqrt{d} \sum_{k=0}^{T-t-1} \rho^k \\ &= \|F\| \delta_s \alpha \rho^{T-t} + \bar{X}_2 \delta_s \alpha \sqrt{d} \left(\frac{1 - \rho^{T-t}}{1 - \rho} \right) \\ &\leq \|F\| \delta_s \alpha \rho^{T-t} + \frac{\bar{X}_2 \delta_s \alpha \sqrt{d}}{1 - \rho}. \end{aligned} \quad (76)$$

Now we bound \bar{q}_t for $t \in \{0, \dots, T\}$. From eq. (65), we have

$$\bar{q}_t = \Phi_t \bar{q}_{t+1} + v_t, \quad (77)$$

where,

$$v_t = \mathbb{E}_{\mu_t} [s_t^\top Y_1(t) s_t] + \begin{bmatrix} \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{1,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \\ \vdots \\ \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{d,t}} [\phi(s_t) \phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \end{bmatrix} + 2\mathbb{E}_{\mu_t} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)] \bar{h}_{t+1}. \quad (78)$$

Given \bar{q}_T , we can write for $t \in \{0, \dots, T-1\}$

$$\bar{q}_t = \prod_{i=T-1}^t \Phi_i \bar{q}_T + \sum_{i=T-1}^t \prod_{j=i-1}^t \Phi_j v_i. \quad (79)$$

Then, we can upper bound $\|\bar{q}_t\|$ as

$$\begin{aligned} \|\bar{q}_t\| &\leq \prod_{i=T-1}^t \|\Phi_i\| \|\bar{q}_T\| + \sum_{i=T-1}^t \prod_{j=i-1}^t \|\Phi_j\| \|v_i\| \\ &\stackrel{(d)}{\leq} \|\bar{q}_T\| + \sum_{i=T-1}^t \|v_i\|, \end{aligned} \quad (80)$$

where in step (d) we have used eq. (72). Now we bound each term of v_t for $t \in \{0, \dots, T-1\}$. We start with

$$\begin{aligned} \|\mathbb{E}_{\mu_{i,t}} [s^\top(t) Y_1(t) s_t]\| &\leq \mathbb{E}_{\mu_{i,t}} [\|s^\top(t) Y_1(t) s_t\|] \\ &\leq \mathbb{E}_{\mu_{i,t}} [\|s_t\|^2 \|Y_1(t)\| \|s_t\|] \\ &\leq \delta_s^2 \|Y_1(t)\|, \end{aligned} \quad (81)$$

for $i \in \{1, \dots, d\}$. Then, we have

$$\|\mathbb{E}_{\mu_t} [s^\top(t)Y_1(t)s_t]\| = \sqrt{\sum_{i=1}^d \|\mathbb{E}_{\mu_{i,t}} [s^\top(t)Y_1(t)s_t]\|^2} \leq \delta_s^2 \|Y_1(t)\| \sqrt{d}. \quad (82)$$

Next, for $i \in \{1, \dots, d\}$, we have

$$\|\bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{i,t}} [\phi(s_t)\phi^\top(s_t)] \otimes Y_2(t)) \bar{h}_{t+1}\| \leq \|\bar{h}_{t+1}\|^2 \|\phi(s_t)\|^2 \|Y_2(t)\| \leq \frac{\|\bar{h}_{t+1}\|^2 \|Y_2(t)\|}{d}. \quad (83)$$

Then,

$$\begin{aligned} \left\| \begin{bmatrix} \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{1,t}} [\phi(s_t)\phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \\ \vdots \\ \bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{d,t}} [\phi(s_t)\phi(s_t)^\top] \otimes Y_2(t)) \bar{h}_{t+1} \end{bmatrix} \right\| &= \sqrt{\sum_{i=1}^d \|\bar{h}_{t+1}^\top (\mathbb{E}_{\mu_{i,t}} [\phi(s_t)\phi^\top(s_t)] \otimes Y_2(t)) \bar{h}_{t+1}\|^2} \\ &\leq \frac{\|\bar{h}_{t+1}\|^2 \|Y_2(t)\|}{\sqrt{d}}. \end{aligned} \quad (84)$$

Next, for $i \in \{1, \dots, d\}$, we have

$$\begin{aligned} \|\mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)]\| &\stackrel{(e)}{\leq} \mathbb{E}_{\mu_{i,t}} [\|\phi(s_t)^\top \otimes s_t^\top Y_3(t)\|] \\ &\leq \mathbb{E}_{\mu_{i,t}} [\|\phi(s_t)\| \|s_t\| \|Y_3(t)\|] \\ &\leq \frac{\delta_s \|Y_3(t)\|}{\sqrt{d}}, \end{aligned} \quad (85)$$

where in step (e) we have used Jensen's inequality. Then,

$$\|\mathbb{E}_{\mu_t} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)]\| = \sqrt{\sum_{i=1}^d \|\mathbb{E}_{\mu_{i,t}} [\phi(s_t)^\top \otimes s_t^\top Y_3(t)]\|^2} \leq \delta_s \|Y_3(t)\|. \quad (86)$$

Let $\|Y_1(t)\| \leq \bar{Y}_1$, $\|Y_2(t)\| \leq \bar{Y}_2$, and $\|Y_3(t)\| \leq \bar{Y}_3$ for $t \in \{0, \dots, T\}$. Using eq. (82), eq. (84), and eq. (86)

$$\|v_t\| \leq \delta_s^2 \bar{Y}_1 \sqrt{d} + \frac{\bar{Y}_2 \|\bar{h}_{t+1}\|^2}{\sqrt{d}} + \delta_s \bar{Y}_3 \|\bar{h}_{t+1}\|, \quad (87)$$

for $t \in \{0, \dots, T\}$. From Theorem A.1, we have $\bar{q}_T = \mathbb{E}_\mu [q_T(s_T)] = \mathbb{E}_\mu [s_T^\top M s_T]$. Then, $\|\bar{q}_T\| \leq \delta_s^2 \|M\|$. Then, we can re-write eq. (80) as

$$\|\bar{q}_t\| \leq \delta_s^2 \|M\| + \delta_s^2 \bar{Y}_1 \sqrt{d} + \frac{\bar{Y}_2 \|\bar{h}_{t+1}\|^2}{\sqrt{d}} + \delta_s \bar{Y}_3 \|\bar{h}_{t+1}\|. \quad (88)$$

□

Following the same notation as eq. (9), let the true parameter be denoted by θ^* , which is written as

$$\theta_t^* = \begin{bmatrix} \theta_{1,t}^* & \dots & \theta_{d,t}^* \end{bmatrix}^\top, \quad \text{where } \theta_{i,t}^* = \begin{bmatrix} \bar{h}_{i,t} \\ \bar{q}_{i,t} \end{bmatrix}. \quad (89)$$

where $\bar{h}_{i,t} \in \mathbb{R}^n$ and $\bar{q}_{i,t} \in \mathbb{R}$ are the components of \bar{h}_t and \bar{q}_t in Theorem C.1 for $i \in \{1, \dots, d\}$ and $t \in \{0, \dots, T\}$.

Corollary C.3. (bound on θ_t^*) Let θ_t^* be as in eq. (89). Then, under the same assumptions of Theorem C.1 and Theorem C.2, we have $\|\theta_t^*\| \leq c_\theta \sqrt{d}$ for $t \in \{0, \dots, T\}$, where $c_\theta > 0$ is independent of d .

Proof. Theorem C.2 implies that for $t \in \{0, \dots, T\}$,

$$\|\bar{h}_t\| \leq a_h + b_h \sqrt{d}, \quad \|\bar{q}_t\| \leq a_q + b_q \sqrt{d}, \quad (90)$$

where $a_h > 0$, $b_h > 0$, $a_q > 0$, and $b_q > 0$ are independent of d . Then, we can bound θ_t^* in eq. (89) as

$$\begin{aligned} \|\theta_t^*\| &= \sqrt{\sum_{i=1}^d (\theta_{i,t}^*)^\top \theta_{i,t}^*} = \sqrt{\sum_{i=1}^d (\bar{h}_{i,t})^\top \bar{h}_{i,t} + (\bar{q}_{i,t})^\top \bar{q}_{i,t}} \\ &= \sqrt{(\bar{h}_t)^\top \bar{h}_t + (\bar{q}_t)^\top \bar{q}_t} = \sqrt{\|\bar{h}_t\|^2 + \|\bar{q}_t\|^2} \\ &\leq \|\bar{h}_t\| + \|\bar{q}_t\| \leq \underbrace{(a_h + b_h + a_q + b_q)}_{c_\theta} \sqrt{d}. \end{aligned} \quad (91)$$

□

Corollary C.3 implies that choosing the projection radius in Algorithm 1 as $R_\theta \geq c_\theta \sqrt{d}$ guarantees that θ_t^* belongs to the projection ball for all t .

D Proof of Theorem 3.3

Let $A_c(t) = A + BK_{x,t}$ and let $\varphi(t_2, t_1) = \prod_{i=t_1}^{t_2-1} A_c(i)$ denote the state transition matrix from t_1 to t_2 .² Let $\pi_t^\ell(x_t, s_t) = K_{x,t}x_t^\ell + K_{s,t}s_t^\ell + K_{h,t}(\phi(s_t)^\top \otimes Z)\theta_{t+1}^\ell$ denote the policy learned from Algorithm 1 at episode ℓ and time t , where $Z = [I_n, 0_{n \times 1}]$. Then, the evolution of x_t in system eq. (1) under the policy $\{\pi_1^\ell, \dots, \pi_{t-1}^\ell\}$ for $t \in \{0, \dots, T\}$ is written as

$$x_t^\ell = \varphi(t, 0)x_0^\ell + \sum_{i=0}^{t-1} \varphi(t, i+1)B\bar{u}_i^\ell, \quad (92)$$

where x_0^ℓ is the initial state at episode ℓ and $\bar{u}_i^\ell = K_{i,s}s_i^\ell(i) + K_{i,h}(\phi(s_i^\ell) \otimes Z)\theta_{i+1}^\ell$. Then,

$$\|x_t^\ell\| \leq \|\varphi(t, 0)\| \|x_0^\ell\| + \|B\| \sum_{i=0}^{t-1} \|\varphi(t, i+1)\| \left(\sup_{0 \leq j \leq t-1} \|\bar{u}_j^\ell\| \right). \quad (93)$$

From (Celi et al., 2022, Lemma B.1), we have $\|\varphi(t_2, t_1)\| \leq \alpha \rho^{t_2-t_1}$ where $\alpha > 0$ and $\rho \in (0, 1)$ are constants. Then, we can write eq. (104) as

$$\begin{aligned} \|x_t^\ell\| &\leq \alpha \rho^t \|x_0^\ell\| + \alpha \|B\| \sum_{i=0}^{t-1} \rho^{t-i-1} \underbrace{\left(\sup_{0 \leq j \leq t-1} \|\bar{u}_j^\ell\| \right)}_{u_\infty^\ell} \\ &\stackrel{(a)}{=} \alpha \rho^t \|x_0^\ell\| + \alpha \|B\| \sum_{k=0}^{t-1} \rho^k u_\infty^\ell \\ &= \alpha \rho^t \|x_0^\ell\| + \alpha \|B\| \left(\frac{1 - \rho^{t-1}}{1 - \rho} \right) u_\infty^\ell \\ &\leq \alpha \rho^t \|x_0^\ell\| + \alpha \|B\| \left(\frac{1}{1 - \rho} \right) u_\infty^\ell, \end{aligned} \quad (94)$$

²The matrix multiplication is performed from the left, i.e., $A(t_1)$ appears as the rightmost matrix in the product.

where in step (a), we have changed the index in the sum to $k = t - i - 1$. Next, we bound on u_∞^ℓ . Let $\|\theta_t^\ell\| \leq R_\theta$, $\|K_{s,t}\| \leq \bar{K}_s$, and $\|K_{h,t}\| \leq \bar{K}_h$ for $t \in \{0, \dots, T-1\}$ and episode ℓ . Then we have,

$$\begin{aligned} \|\bar{u}_t^\ell\| &\leq \|K_{s,t}\| \|s_t^\ell\| + \|K_{h,t}\| (\phi(s_t^\ell) \otimes Z) \|\theta_{t+1}^\ell\| \\ &\leq \bar{K}_s \delta_s + \bar{K}_h \|\phi(s_t^\ell)\| \|Z\| R_\theta \\ &\leq \bar{K}_s \delta_s + \frac{\bar{K}_h R_\theta}{\sqrt{d}}. \end{aligned} \tag{95}$$

Since the above bound is uniform for $t \in \{0, \dots, T-1\}$, we have $u_\infty^\ell \leq \bar{K}_s \delta_s + \frac{\bar{K}_h R_\theta}{\sqrt{d}}$. The proof follows by substituting the bound of u_∞^ℓ in eq. (94).

E Proof of Theorem 3.4

We begin by presenting the following technical Lemmas.

Lemma E.1. *Let $X_t = \sum_{i=1}^{t-1} z_i z_i^\top + \gamma I_p$, where $z_i \in \mathbb{R}^p$ and $\gamma > 0$. Let $\mathbb{E}[zz^\top] \succeq \alpha I_p$ with $\alpha > 0$, and $\|z\| \leq \zeta$. Let $\delta \in [0, 1]$ and assume $t \geq (8\zeta^2 \log(p/\delta))/\alpha$. Then, with probability at least $1 - \delta$, the minimum eigenvalue of X_t satisfies*

$$\lambda_{\min}(X_t) \geq \gamma + \frac{(t-1)\alpha}{2}.$$

Proof. Let $\alpha = \lambda_{\min}(\mathbb{E}[zz^\top])$. Define

$$\mu_{\min} \triangleq \lambda_{\min} \left(\sum_{i=1}^{t-1} \mathbb{E}[zz^\top] \right) = \lambda_{\min}((t-1)\mathbb{E}[zz^\top]) = (t-1)\lambda_{\min}(\mathbb{E}[zz^\top]) = (t-1)\alpha.$$

Further, we have $z_i z_i^\top \succeq 0$ and $\lambda_{\max}(z_i z_i^\top) = \|z_i\|^2 \leq \zeta^2$. Then, using (Tropp, 2012, Theorem 1.1), we have

$$\begin{aligned} \mathbb{P} \left(\lambda_{\min} \left(\sum_{i=1}^{t-1} z_i z_i^\top \right) \leq (1-\varepsilon)(t-1)\alpha \right) &\leq p \left(\frac{\exp(-\varepsilon)}{(1-\varepsilon)^{1-\varepsilon}} \right)^{\frac{(t-1)\alpha}{\zeta^2}} \\ &\stackrel{(a)}{\leq} p \exp \left(\frac{-\varepsilon^2(t-1)\alpha}{2\zeta^2} \right), \end{aligned} \tag{96}$$

for $\varepsilon \in [0, 1]$, where in step (a) we have used $\frac{\exp(-\varepsilon)}{(1-\varepsilon)^{1-\varepsilon}} \leq \exp(-\varepsilon^2/2)$ for $\varepsilon \in (0, 1)$. Choose $\varepsilon = 0.5$, then we write eq. (96) as

$$\mathbb{P} \left(\lambda_{\min} \left(\sum_{i=1}^{t-1} z_i z_i^\top \right) \leq \frac{(t-1)\alpha}{2} \right) \leq p \exp \left(\frac{-(t-1)\alpha}{8\zeta^2} \right). \tag{97}$$

Let $p \exp \left(\frac{-(t-1)\alpha}{8\zeta^2} \right) \leq \delta$, then we have

$$t \geq \frac{8\zeta^2 \log(p/\delta)}{\alpha}.$$

Then, with probability at least $1 - \delta$ we have

$$\lambda_{\min} \left(\sum_{i=1}^{t-1} z_i z_i^\top \right) \geq \frac{(t-1)\alpha}{2}.$$

Finally, we have

$$\lambda_{\min}(X_t) \geq \lambda_{\min} \left(\sum_{i=1}^{t-1} z_i z_i^\top \right) + \gamma \geq \frac{(t-1)\alpha}{2} + \gamma.$$

□

Lemma E.2. Consider the system eq. (1) and the Markov process eq. (2). Let Assumption 2.1 be satisfied, and let

$$x_{t+1} = \varphi(t+1, 0)x_0 + \sum_{i=0}^t \varphi(t+1, i+1)B\bar{u}_i(s_i),$$

with $x_0 \sim \mathcal{N}(0, \Sigma_0)$, and $\bar{u}_i(s_i)$ is an arbitrary input that depends on s_i and is independent of x_0 . Let

$$\psi_t = \phi(s_t) \otimes \begin{bmatrix} 2x_{t+1} \\ 1 \end{bmatrix}.$$

Assume $\Sigma_0 \succ 0$ and $\varphi(t+1, 0)$ is nonsingular for $t \in \{0, \dots, T-1\}$. Then, $\mathbb{E}[\psi_t \psi_t^\top] \succ 0$ for $t \in \{0, \dots, T-1\}$.

Proof. We begin by writing

$$\psi_t = \phi(s_t) \otimes \begin{bmatrix} 2x_{t+1} \\ 1 \end{bmatrix} = \begin{bmatrix} 2\phi(s_t) \otimes x_{t+1} \\ \phi(s_t) \end{bmatrix}.$$

Then,

$$\psi_t \psi_t^\top = \begin{bmatrix} 4\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}x_{t+1}^\top & 2\phi(s_t)\phi(s_t)^\top \otimes x_{t+1} \\ 2\phi(s_t)\phi(s_t)^\top \otimes x_{t+1} & \phi(s_t)\phi(s_t)^\top \end{bmatrix}.$$

Taking the expectation, we get

$$\mathbb{E}[\psi_t \psi_t^\top] = \begin{bmatrix} 4\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}x_{t+1}^\top] & 2\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}] \\ 2\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}] & \mathbb{E}[\phi(s_t)\phi(s_t)^\top] \end{bmatrix}.$$

From Assumption 2.1, we have $\mathbb{E}[\phi(s_t)\phi(s_t)^\top] \succ 0$ for all t . For notational convenience we denote $\Sigma_\phi = \mathbb{E}[\phi(s_t)\phi(s_t)^\top]$. We apply the Schur complement

$$S = 4\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}x_{t+1}^\top] - 4\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}] \Sigma_\phi^{-1} \mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}]. \quad (98)$$

Showing $\mathbb{E}[\psi_t \psi_t^\top] \succ 0$ boils down to showing that $S \succ 0$. From eq. (98), we have

$$\begin{aligned} \mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}x_{t+1}^\top] &\stackrel{(a)}{=} \mathbb{E}[\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes x_{t+1}x_{t+1}^\top | s_t]] \\ &= \mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes \mathbb{E}[x_{t+1}x_{t+1}^\top | s_t]] \\ &= \mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes \Sigma_{x|s}] + \mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes \mu_x(s_t)\mu_x(s_t)^\top], \end{aligned} \quad (99)$$

where in step (a) we used the law of total expectation, and

$$\begin{aligned} \Sigma_{x|s} &= \mathbb{E}[(x_{t+1} - \mathbb{E}[x_{t+1}|s_t])(x_{t+1} - \mathbb{E}[x_{t+1}|s_t])^\top | s_t], \\ \mu_x(s_t) &= \mathbb{E}[x_{t+1}|s_t]. \end{aligned}$$

For notational convenience, let $z = \phi(s_t) \otimes \mu_x(s_t)$. Substituting eq. (99) in eq. (98), we get

$$S = \underbrace{4\mathbb{E}[\phi(s_t)\phi(s_t)^\top \otimes \Sigma_{x|s}]}_{S_1} + \underbrace{4\mathbb{E}[zz^\top] - 4\mathbb{E}[z\phi(s_t)^\top] \Sigma_\phi^{-1} \mathbb{E}[\phi(s_t)z^\top]}_{S_2}. \quad (100)$$

We have $S_1 \succeq 0$ since $\phi(s_t)\phi(s_t)^\top \succeq 0$ and $\Sigma_{x|s} \succeq 0$. From eq. (100), we have

$$\begin{aligned}
S_2 &= 4\mathbb{E} [zz^\top] - 4\mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] \\
&\stackrel{(b)}{=} 4\mathbb{E} [zz^\top] - 4\mathbb{E} \left[z\phi(s_t)^\top \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] \right] + 4\mathbb{E} \left[\mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t)z^\top \right] \\
&\quad - 4\mathbb{E} \left[\mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t)z^\top \right] \\
&\stackrel{(c)}{=} 4\mathbb{E} [zz^\top] - 4\mathbb{E} \left[z\phi(s_t)^\top \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] \right] + 4\mathbb{E} \left[z\phi(s_t)^\top \Sigma_\phi^{-1} \Sigma_\phi \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] \right] \\
&\quad - 4\mathbb{E} \left[\mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t)z^\top \right] \\
&= 4\mathbb{E} \left[zz^\top - z\phi(s_t)^\top \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] - \mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t)z^\top \right] \\
&\quad + \mathbb{E} \left[z\phi(s_t)^\top \Sigma_\phi^{-1} \phi(s_t)\phi(s_t)^\top \Sigma_\phi^{-1} \mathbb{E} \left[\phi(s_t)z^\top \right] \right] \\
&= 4\mathbb{E} \left[\left(z - \mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t) \right) \left(z - \mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t) \right)^\top \right],
\end{aligned}$$

where in step (b) we have added and subtracted the term $4\mathbb{E} \left[\mathbb{E} \left[z\phi(s_t)^\top \right] \Sigma_\phi^{-1} \phi(s_t)z^\top \right]$, and in step (c) we have used $I = \Sigma_\phi \Sigma_\phi^{-1}$. Then, we have $S = S_1 + S_2 \succeq 0$. From eq. (92) we have $x_{t+1} = \varphi(t+1, 0)x_0 + \sum_{i=0}^t \varphi(t+1, i+1)B\bar{u}_i$, hence, $\mathbb{E} [x_{t+1}|s_t] = \mathbb{E} \left[\sum_{i=0}^t \varphi(t+1, i+1)B\bar{u}_i | s_t \right]$. For notational convenience, let $\tilde{u}(t) = \sum_{i=0}^t \varphi(t+1, i+1)B\bar{u}_i$. Then we get

$$\begin{aligned}
\Sigma_{x|s} &= \varphi(t+1, 0)\Sigma_0\varphi(t+1, 0)^\top + \mathbb{E} \left[(\tilde{u}(t) - \mathbb{E} [\tilde{u}(t)|s_t]) (\tilde{u}(t) - \mathbb{E} [\tilde{u}(t)|s_t])^\top | s_t \right] \\
&\succeq \varphi(t+1, 0)\Sigma_0\varphi(t+1, 0)^\top.
\end{aligned}$$

Since $\Sigma_0 \succ 0$ and $\varphi(t+1, 0)$ is nonsingular for all t , then, $\varphi(t+1, 0)\Sigma_0\varphi(t+1, 0)^\top \succ 0$. Then, we have $S_1 = 4\mathbb{E} \left[\phi(s_t)\phi(s_t)^\top \otimes \Sigma_{x|s} \right] \succeq 4\mathbb{E} \left[\phi(s_t)\phi(s_t)^\top \right] \otimes \left(\varphi(t+1, 0)\Sigma_0\varphi(t+1, 0)^\top \right) \succ 0$ since $\mathbb{E} \left[\phi(s_t)\phi(s_t)^\top \right] \succ 0$, and $\varphi(t+1, 0)\Sigma_0\varphi(t+1, 0)^\top$ is independent of s . Therefore, $S \succ 0$, which implies $\mathbb{E} \left[\psi_t \psi_t^\top \right] \succ 0$ for $t \in \{0, \dots, T-1\}$. \square

Now we present the proof of Theorem 3.4. For notational convenience, we denote $\phi(s_t^i)$ and $Y(x_t^i, u_t^i)$ by ϕ_t^i and Y_t^i , respectively. We have from the expression of θ_{t+1}^ℓ in eq. (15)

$$\begin{aligned}
\epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i) &= 2x_{t+1}^i h_{t+1}^\ell(s_{t+1}^i) + q_{t+1}^\ell(s_{t+1}^i) \\
&= \underbrace{\begin{bmatrix} 2x_{t+1}^i & 1 \end{bmatrix}}_{(y_t^i)^\top} \underbrace{\begin{bmatrix} h_{t+1}^\ell(s_{t+1}^i) \\ q_{t+1}^\ell(s_{t+1}^i) \end{bmatrix}}_{v_{t+1}^\ell(s_{t+1}^i)}.
\end{aligned}$$

We can derive an upper bound on $|\epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i)|$ as

$$|\epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i)| \leq 2\|x_{t+1}^i\| \|h_{t+1}^\ell(s_{t+1}^i)\| + \|q_{t+1}^\ell(s_{t+1}^i)\|. \quad (101)$$

Using eq. (49) we can write

$$\begin{aligned}
\|h_{t+1}^\ell(s_{t+1}^i)\| &\leq \|X_1(t+1)\| \|\phi_{t+1}^i\| \|\theta_{t+2}^\ell\| + \|X_2(t+1)\| \|s_{t+1}^i\| \\
&\leq \frac{\bar{X}_1 R_\theta}{\sqrt{d}} + \bar{X}_2 \delta_s,
\end{aligned} \quad (102)$$

where $\|X_1(t)\| \leq \bar{X}_1$ and $\|X_2(t)\| \leq \bar{X}_2$ for all $t \in \{0, \dots, T\}$. From eq. (50) we can write

$$\begin{aligned} q_{t+1}^\ell(s_{t+1}^i) &\leq \|\phi_{t+1}^i\| \|\theta_{t+2}^\ell\| + \|s_{t+1}^i\|^2 \|Y_1(t+1)\| \\ &\quad + \|\theta_{t+2}^\ell\|^2 \|\phi_{t+1}^i\|^2 \|Y_2(t+1)\| \\ &\quad + 2\|s_{t+1}^i\| \|Y_3(t+1)\| \|\phi_{t+1}^i\| \|\theta_{t+2}^\ell\| \\ &\leq \frac{R_\theta}{\sqrt{d}} + \delta_s^2 \bar{Y}_1 + \frac{R_\theta^2 \bar{Y}_2}{d} + 2 \frac{\delta_s \bar{Y}_3 R_\theta}{\sqrt{d}}, \end{aligned} \quad (103)$$

where $\|Y_1(t)\| \leq \bar{Y}_1$, $\|Y_2(t)\| \leq \bar{Y}_2$, and $\|Y_3(t)\| \leq \bar{Y}_3$ for all $t \in \{0, \dots, T\}$. From Theorem 3.3, we have

$$\|x_t^\ell\| \leq \alpha \rho^t \|x^\ell(0)\| + \frac{\alpha \|B\|}{1-\rho} \left(\bar{K}_s \delta_s + \frac{\bar{K}_h R_\theta}{\sqrt{d}} \right),$$

for $t \in \{0, \dots, T\}$ and $\ell \in \{1, \dots, L\}$, with $\alpha > 0$ and $0 < \rho < 1$. Define

$$\bar{x} = \sup_{\substack{t \in \{0, \dots, T\} \\ \ell \in \{1, \dots, L\}}} \left\{ \alpha \rho^t \|x^\ell(0)\| + \frac{\alpha \|B\|}{1-\rho} \left(\bar{K}_s \delta_s + \frac{\bar{K}_h R_\theta}{\sqrt{d}} \right) \right\}. \quad (104)$$

Substituting eq. (102), eq. (103), and eq. (104) in eq. (101), we get

$$|\epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i)| \leq \left(\frac{2\bar{x}\bar{X}_1 + 1 + 2\delta_s \bar{Y}_3}{\sqrt{d}} \right) R_\theta + \frac{\bar{Y}_2}{d} R_\theta^2 + 2\bar{X}_2 \bar{x} \delta_s + \bar{Y}_1 \delta_s^2. \quad (105)$$

Further, we can bound

$$\begin{aligned} \|\psi_t^i\| &= \|(Y_t^i)^\top \phi_t^i\| \leq \|[2x_{t+1}^i \quad 1]\| \|\phi_t^i\| \\ &\leq \sqrt{\frac{4\bar{x}^2 + 1}{d}} \triangleq \delta_\psi \end{aligned} \quad (106)$$

Next, using the expression of θ_{t+1}^ℓ in eq. (15), we write

$$\theta_{t+1}^\ell - \theta_{t+1}^* = (\Lambda_t^\ell)^{-1} \sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i \epsilon_{t+1}^\ell(x_{t+1}^i, s_{t+1}^i) - \theta_{t+1}^*, \quad (107)$$

We re-write eq. (107) as

$$\begin{aligned} \theta_{t+1}^\ell - \theta_{t+1}^* &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} (Y_t^i)^\top (\phi_t^i) (y_t^i)^\top v_{t+1}^\ell(s_{t+1}^i) - \Lambda_t^\ell \theta_{t+1}^* \right) \\ &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} (Y_t^i)^\top (\phi_t^i) (y_t^i)^\top v_{t+1}^\ell(s_{t+1}^i) - \sum_{i=1}^{\ell-1} (Y_t^i)^\top (\phi_t^i) (\phi_t^i)^\top Y_t^i \theta_{t+1}^* \right) - \lambda (\Lambda_t^\ell)^{-1} \theta_{t+1}^*. \end{aligned} \quad (108)$$

From eq. (108), we expand the term as

$$\begin{aligned} (\phi_t^i)^\top Y_t^i \theta_{t+1}^* &= \sum_{j=1}^d \phi_{j,t}^i (y_t^i)^\top \begin{bmatrix} h_{j,t+1}^* \\ q_{j,t+1}^* \end{bmatrix} \\ &= (y_t^i)^\top \sum_{j=1}^d \phi_{j,t}^i \mathbb{E}_{\mu_j} \left[\begin{bmatrix} h_{t+1}^*(s_{t+1}) \\ q_{t+1}^*(s_{t+1}) \end{bmatrix} \right] \\ &= (y_t^i)^\top \sum_{j=1}^d \phi_{j,t}^i \int_{\mathcal{S}} \underbrace{\begin{bmatrix} h_{t+1}^*(s_{t+1}) \\ q_{t+1}^*(s_{t+1}) \end{bmatrix}}_{v_{t+1}^*(s_{t+1})} \mu_j(ds_{t+1}) \\ &= (y_t^i)^\top \int_{\mathcal{S}} v_{t+1}^*(s_{t+1}) \sum_{j=1}^d \phi_{j,t}^i \mu_j(ds_{t+1}) \\ &= (y_t^i)^\top \mathbb{E} [v_{t+1}^*(s_{t+1}) | s_t^i]. \end{aligned} \quad (109)$$

For notational convenience, we use $\psi_t^i = (Y_t^i)^\top \phi_t^i$. Substituting eq. (109) in eq. (108), we get

$$\begin{aligned} \theta_{t+1}^\ell - \theta_{t+1}^* &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} \psi_t^i (y_t^i)^\top \left(v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E} [v_{t+1}^*(s_{t+1}) | s_t^i] \right) \right) - \lambda (\Lambda_t^\ell)^{-1} \theta_{t+1}^* \\ &= (\Lambda_t^\ell)^{-1} \underbrace{\left(\sum_{i=1}^{\ell-1} \psi_t^i (y_t^i)^\top \left(v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E} [v_{t+1}^\ell(s_{t+1}) | s_t^i] \right) \right)}_{r_1} \\ &\quad + (\Lambda_t^\ell)^{-1} \underbrace{\left(\sum_{i=1}^{\ell-1} \psi_t^i (y_t^i)^\top \left(\mathbb{E} [v_{t+1}^\ell(s_{t+1}) - v_{t+1}^*(s_{t+1}) | s_t^i] \right) \right)}_{r_2} - \underbrace{\lambda (\Lambda_t^\ell)^{-1} \theta_{t+1}^*}_{r_3}. \end{aligned} \quad (110)$$

Let $\xi_{s,t+1}^\ell = (y_t^i)^\top v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E} [(y_t^i)^\top v_{t+1}^\ell(s_{t+1}) | s_t^i]$. We have $\mathbb{E} [\xi_{s,t+1}^\ell | \mathcal{F}_t^{\ell-1}] = 0$. Since $(y_t^i)^\top v_{t+1}^\ell(s_{t+1}^i)$ is bounded (see eq. (105)), then ξ is σ -subGaussian with

$$\sigma = \left(\frac{2\bar{x}\bar{X}_1 + 1 + 2\delta_s\bar{Y}_3}{\sqrt{d}} \right) R_\theta + \frac{\bar{Y}_2}{d} R_\theta^2 + 2\bar{X}_2\bar{x}\delta_s + \bar{Y}_1\delta_s^2. \quad (111)$$

Then, using (Abbasi-Yadkori et al., 2011, Theorem 1), we have with probability at least $1 - \delta$, we have

$$\left\| \sum_{i=1}^{\ell-1} (\psi_t^i (y_t^i)^\top (v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E} [v_{t+1}^\ell(s_{t+1}) | s_t^i])) \right\|_{(\Lambda_t^\ell)^{-1}}^2 \leq 2\sigma^2 \left(\log \left(\sqrt{\frac{\det(\Lambda_t^\ell)}{\det(\Lambda_t^1)}} \right) + \log \left(\frac{1}{\delta} \right) \right). \quad (112)$$

Recall from Alg. 1, we have

$$\begin{aligned} \Lambda_t^1 &= \lambda I_{d(n+1)}, \\ \Lambda_t^\ell &= \sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i (\phi_t^i)^\top Y_t^i + \lambda I_{d(n+1)} = \lambda \underbrace{\left(\frac{1}{\lambda} \sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i (\phi_t^i)^\top Y_t^i + I_{d(n+1)} \right)}_{\bar{\Lambda}_t^\ell}. \end{aligned}$$

Then,

$$\frac{\det(\Lambda_t^\ell)}{\det(\Lambda_t^1)} = \frac{\det(\lambda \bar{\Lambda}_t^\ell)}{\det(\lambda I_{d(n+1)})} = \det(\bar{\Lambda}_t^\ell) = \det \left(\frac{1}{\lambda} \sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i (\phi_t^i)^\top Y_t^i + I_{d(n+1)} \right) = \prod_{i=1}^{d(n+1)} (1 + \gamma_i), \quad (113)$$

where γ_i is the i -th eigenvalue of $\frac{1}{\lambda} \sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i (\phi_t^i)^\top Y_t^i$. Let $\|(Y_t^i)^\top \phi_t^i\| \leq \delta_\psi$ for $t \in \{0, \dots, T\}$ and $i \in \{1, \dots, L\}$ (see eq. (106)). Then,

$$\gamma_i \leq \frac{1}{\lambda} \sum_{i=1}^{\ell-1} \text{Tr} (Y_t^i)^\top \phi_t^i (\phi_t^i)^\top Y_t^i \leq \frac{1}{\lambda} \sum_{i=1}^{\ell-1} \|(Y_t^i)^\top \phi_t^i\|^2 \leq \frac{\ell \delta_\psi^2}{\lambda}. \quad (114)$$

Then

$$\log \left(\sqrt{\det(\bar{\Lambda}_t^\ell)} \right) = \frac{1}{2} \log \left(\prod_{i=1}^{d(n+1)} (1 + \gamma_i) \right) = \frac{1}{2} \sum_{i=1}^{d(n+1)} \log(1 + \gamma_i) \leq \frac{d(n+1)}{2} \log \left(1 + \frac{\ell \delta_\psi^2}{\lambda} \right). \quad (115)$$

Substituting eq. (115) in eq. (112), we get with probability at least $1 - \delta$

$$\left\| \sum_{i=1}^{\ell-1} (\psi_t^i (y_t^i)^\top (v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E} [v_{t+1}^\ell(s_{t+1}) | s_t^i])) \right\|_{(\Lambda_t^\ell)^{-1}}^2 \leq \sigma^2 \left(d(n+1) \log \left(1 + \frac{\ell \delta_\psi^2}{\lambda} \right) + 2 \log \left(\frac{1}{\delta} \right) \right). \quad (116)$$

Then,

$$\left\| \sum_{i=1}^{\ell-1} (\psi_t^i(y_t^i))^\top (v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E}[v_{t+1}^\ell(s_{t+1}^i)|s_t^i]) \right\|_{(\Lambda_t^\ell)^{-1}} \leq \sigma \sqrt{d(n+1) \log \left(1 + \frac{\ell \delta_\psi^2}{\lambda} \right) + 2 \log \left(\frac{1}{\delta} \right)}. \quad (117)$$

Then, we have

$$\begin{aligned} |\phi_t^\top Y_t r_1| &\leq \left\| \phi_t^\top Y_t (\Lambda_t^\ell)^{-\frac{1}{2}} \right\| \left\| \sum_{i=1}^{\ell-1} (\psi_t^i(y_t^i))^\top (v_{t+1}^\ell(s_{t+1}^i) - \mathbb{E}[v_{t+1}^\ell(s_{t+1}^i)|s_t^i]) \right\|_{(\Lambda_t^\ell)^{-1}} \\ &\leq \sigma \sqrt{\left(d(n+1) \log \left(1 + \frac{\ell \delta_\psi^2}{\lambda} \right) + 2 \log \left(\frac{1}{\delta} \right) \right)} \sqrt{\phi_t^\top Y_t (\Lambda_t^\ell)^{-1} Y_t^\top \phi_t}. \end{aligned} \quad (118)$$

Next, we have

$$\begin{aligned} |\phi_t^\top Y_t r_3| &\leq \lambda \left\| \phi_t^\top Y_t (\Lambda_t^\ell)^{-\frac{1}{2}} \right\| \left\| (\Lambda_t^\ell)^{-\frac{1}{2}} \right\| \|\theta_{t+1}^*\| \\ &\leq \sqrt{\lambda} \|\theta_{t+1}^*\| \sqrt{\phi_t^\top Y_t (\Lambda_t^\ell)^{-1} Y_t^\top \phi_t}. \end{aligned} \quad (119)$$

Next, we have

$$\begin{aligned} r_2 &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i(y_t^i)^\top (\mathbb{E}[v_{t+1}^\ell(s_{t+1}^i)|s_t^i] - \mathbb{E}[v_{t+1}^*(s_{t+1}^i)|s_t^i]) \right) \\ &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i(y_t^i)^\top \int_{\mathcal{S}} (v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)) \sum_{j=1}^d \phi_{j,t}^i \mu_{j,t}(ds_{t+1}^i) \right) \\ &= (\Lambda_t^\ell)^{-1} \left(\sum_{i=1}^{\ell-1} (Y_t^i)^\top \phi_t^i(y_t^i)^\top \sum_{j=1}^d \phi_{j,t}^i \int_{\mathcal{S}} (v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)) \mu_{j,t}(ds_{t+1}^i) \right) \\ &= (\Lambda_t^\ell)^{-1} \underbrace{\left(\sum_{i=1}^{\ell-1} \psi_t^i(\psi_t^i)^\top \right)}_{\Lambda_t^\ell - \lambda I_{d(n+1)}} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)] \\ &= \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)] - \lambda (\Lambda_t^\ell)^{-1} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)]. \end{aligned} \quad (120)$$

Then, we have

$$\begin{aligned} \phi_t^\top Y_t r_2 &= \phi_t^\top Y_t \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)] - \lambda \phi_t^\top Y_t (\Lambda_t^\ell)^{-1} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)] \\ &= y_t^\top \mathbb{E} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i) | s_t] - \lambda \phi_t^\top Y_t (\Lambda_t^\ell)^{-1} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)] \\ &= \mathbb{E} [V_{t+1}^\ell(x_{t+1}, s_{t+1}) - V_{t+1}(x_{t+1}, s_{t+1}) | s_t] - \underbrace{\lambda \phi_t^\top Y_t (\Lambda_t^\ell)^{-1} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}^i) - v_{t+1}^*(s_{t+1}^i)]}_{r_4}. \end{aligned} \quad (121)$$

Since we choose R_θ in Algorithm 1 to be the upper bound on $\|\theta_{t+1}^*\|$ (which we derive in Appendix C.2) for $t \in \{0, \dots, T-1\}$, and we have $\|\theta_{t+1}^\ell\| \leq R_\theta$ (from Algorithm 1) for $t \in \{0, \dots, T-1\}$ and $\ell \in \{1, \dots, L\}$, it can be seen from eq. (102) and eq. (103) that $\|v_{t+1}^\ell(s_{t+1}^i)\| \leq \delta_v$ and $\|v_{t+1}^*(s_{t+1}^i)\| \leq \delta_v$. We can derive an expression for δ_v using eq. (102) and eq. (103) as

$$\|v_{t+1}^\ell(s_{t+1}^i)\| = \left\| \begin{bmatrix} h_{t+1}^\ell(s_{t+1}^i) \\ q_{t+1}^\ell(s_{t+1}^i) \end{bmatrix} \right\| \leq \underbrace{\sqrt{\|h_{t+1}^\ell(s_{t+1}^i)\|^2 + \|q_{t+1}^\ell(s_{t+1}^i)\|^2}}_{\delta_v}. \quad (122)$$

Then, we can bound $|r_4|$ in eq. (121) as

$$\begin{aligned} |\lambda \phi_t^\top Y_t (\Lambda_t^\ell)^{-1} \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell(s_{t+1}) - v_{t+1}^*(s_{t+1})]| &\leq \lambda \left\| \phi_t^\top Y_t (\Lambda_t^\ell)^{-\frac{1}{2}} \right\| \left\| (\Lambda_t^\ell)^{-\frac{1}{2}} \right\| \left\| \mathbb{E}_{\mu_{t+1}} [v_{t+1}^\ell - v_{t+1}^*] \right\| \\ &\leq 2\sqrt{\lambda} \delta_v \sqrt{\phi_t^\top Y_t (\Lambda_t^\ell)^{-1} Y_t^\top \phi_t}. \end{aligned} \quad (123)$$

Using the parametrized form of the Q -function from Theorem 3.1, we have

$$\begin{aligned} Q_t^\ell(x_t, s_t, u_t) - Q_t^*(x_t, s_t, u_t) &= \phi_t^\top Y_t (\theta_{t+1}^\ell - \theta_{t+1}^*) = \phi_t^\top Y_t (r_1 + r_2 + r_3) \\ \implies Q_t^\ell(x_t, s_t, u_t) - Q_t^*(x_t, s_t, u_t) - \mathbb{E} [V_{t+1}^\ell(x_{t+1}, s_{t+1}) - V_{t+1}^*(x_{t+1}, s_{t+1}) | s_t] &= \underbrace{\phi_t^\top Y_t (r_1 + r_3) - r_4}_{\Delta_t^\ell(x_t, s_t, u_t)}. \end{aligned} \quad (124)$$

Then, using eq. (118), eq. (119), and eq. (123), we can bound

$$\begin{aligned} |\Delta_t^\ell(x_t, s_t, u_t)| &\leq \left(\sigma \sqrt{\left(d(n+1) \log \left(1 + \frac{\ell \delta_\psi^2}{\lambda} \right) + 2 \log \left(\frac{1}{\delta} \right) \right)} + \|\theta_{t+1}^*\| \sqrt{\lambda} + 2\sqrt{\lambda} \delta_v \right) \sqrt{\phi_t^\top Y_t (\Lambda_t^\ell)^{-1} Y_t^\top \phi_t} \\ &= \chi(\ell) \sqrt{\phi_t^\top Y_t (\Lambda_t^\ell)^{-1} Y_t^\top \phi_t}. \end{aligned} \quad (125)$$

Let $\delta_t^\ell = V_t^\ell(x_t^{*\ell}, s_t^\ell) - V_t^*(x_t^{*\ell}, s_t^\ell)$ and $\zeta_{t+1}^\ell = \mathbb{E} [\delta_{t+1}^\ell | s_t^\ell] - \delta_{t+1}^\ell$, where $x_t^{*\ell}$ is the state under the optimal policy π_t^* starting from x_0^ℓ for $t \in \{0, \dots, T\}$ and $\ell \in \{1, \dots, L\}$. From the definition of the value function, we have $V_t^*(x, s) = \min_{u \in \mathcal{U}} Q_t^*(x, s, u)$. Further, since Algorithm 1 selects a greedy policy with respect to $Q_t^\ell(x, s, u)$, we have $V_t^\ell(x, s) = \min_{u \in \mathcal{U}} Q_t^\ell(x, s, u)$. Let $u_t^{*\ell}$ be the optimal control input at that generates $x_t^{*\ell}$ for $t \in \{0, \dots, T-1\}$ and $\ell \in \{1, \dots, L\}$. Then, we can write

$$\begin{aligned} \delta_t^\ell &= V_t^\ell(x_t^{*\ell}, s_t^\ell) - V_t^*(x_t^{*\ell}, s_t^\ell) \\ &= Q_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) - Q_t^*(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) \\ &\leq Q_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) - Q_t^*(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}). \end{aligned} \quad (126)$$

Then, from eq. (124) we can write

$$Q_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) - Q_t^*(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) = \mathbb{E} [\delta_{t+1}^\ell | s_t^\ell] + \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) = \delta_{t+1}^\ell + \zeta_{t+1}^\ell + \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) \quad (127)$$

Substituting eq. (126) in eq. (127), we get

$$\delta_t^\ell \leq \delta_{t+1}^\ell + \zeta_{t+1}^\ell + \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}).$$

Note that $x_0^{*\ell} = x_0^\ell$. Next, from eq. (17) we write

$$\mathcal{R}(L) = \sum_{\ell=1}^L \delta_0^\ell \leq \sum_{\ell=1}^L \sum_{t=1}^T \zeta_t^\ell + \sum_{\ell=1}^L \sum_{t=0}^{T-1} \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}). \quad (128)$$

We have $\mathbb{E} [\zeta_t^\ell | \mathcal{F}_{t-1}^\ell] = 0$. Also, since we have $|\delta_t^\ell| \leq \sigma$ (see eq. (105)) with σ as in eq. (111), then we have $|\zeta_t^\ell| \leq \sigma$ for $t \in \{0, \dots, T\}$ and $\ell \in \{0, \dots, L\}$. Hence ζ_t^ℓ is a martingale difference sequence. Then, using the Azuma-Hoeffding inequality, for any $\varepsilon > 0$, we get

$$\mathbb{P} \left(\sum_{\ell=1}^L \sum_{t=1}^T \zeta_t^\ell \geq \varepsilon \right) \leq \exp \left(\frac{-\varepsilon^2}{2 \sum_{i=1}^{LT} \sigma^2} \right) = \delta.$$

Hence, we get with probability at least $1 - \delta$

$$\sum_{\ell=1}^L \sum_{t=1}^T \zeta_t^\ell \leq \sigma \sqrt{2LT \log(1/\delta)}. \quad (129)$$

Next, using eq. (125) we can write

$$\begin{aligned}
\sum_{\ell=1}^L \sum_{t=0}^{T-1} \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) &\leq \sum_{\ell=1}^L \sum_{t=0}^{T-1} \chi(\ell) \sqrt{(\phi_t^\ell)^\top Y_t^{*\ell} (\Lambda_t^\ell)^{-1} (Y_t^{*\ell})^\top \phi_t(s_t^\ell)} \\
&\leq \chi(L) \sum_{\ell=1}^L \sum_{t=0}^{T-1} \|(\Lambda_t^\ell)^{-1/2} (Y_t^{*\ell})^\top \phi_t(s_t^\ell)\| \\
&\stackrel{(a)}{\leq} \chi(L) \delta_\psi \sum_{\ell=1}^L \sum_{t=0}^{T-1} \|(\Lambda_t^\ell)^{-1/2}\|,
\end{aligned}$$

where in step (a) we have used eq. (106). Assume the state transition matrix, $\varphi(t, 0)$, in eq. (92) is nonsingular for $t \in \{0, \dots, T\}$.³ Then, from Lemma E.2 we have $\mathbb{E}[\psi_t \psi_t^\top] \succeq \gamma I_{d(n+1)}$ with $\gamma > 0$. Further, from eq. (106), we have $\|\psi_t^\ell\| \leq \delta_\psi$ for $t \in \{0, \dots, T\}$ and $\ell \in \{1, \dots, L\}$. Let $\delta \in [0, 1]$ and assume $\ell \geq (8\delta_\psi^2 \log(d(n+1)/\delta))/\gamma$. Then, using Lemma E.1 we have with probability at least $1 - \delta$

$$\lambda_{\min}(\Lambda_t^\ell) \geq \lambda + \frac{(\ell-1)\gamma}{2}.$$

Then, we can write

$$\begin{aligned}
\sum_{\ell=1}^L \sum_{t=0}^{T-1} \|(\Lambda_t^\ell)^{-1/2}\| &\leq \sum_{\ell=1}^L \sum_{t=0}^{T-1} \frac{1}{\sqrt{\lambda_{\min}(\Lambda_t^\ell)}} \\
&\leq \sum_{\ell=1}^L \frac{T}{\sqrt{\lambda + \frac{(\ell-1)\gamma}{2}}} \\
&\leq \frac{T}{\sqrt{\lambda}} + \sum_{\ell=2}^L \frac{T}{\sqrt{\lambda + \frac{(\ell-1)\gamma}{2}}} \\
&\leq \frac{T}{\sqrt{\lambda}} + \sum_{\ell=2}^L \frac{T\sqrt{2}}{\sqrt{(\ell-1)\gamma}} \\
&\leq \frac{T}{\sqrt{\lambda}} + \frac{2\sqrt{2}T}{\sqrt{\gamma}} \sum_{\ell=2}^L (\sqrt{(\ell-1)} - \sqrt{(\ell-2)}) \\
&= \frac{T}{\sqrt{\lambda}} + \frac{2\sqrt{2}T}{\sqrt{\gamma}} \sqrt{L-1} \\
&\leq \frac{T}{\sqrt{\lambda}} + \frac{4T\sqrt{L}}{\sqrt{\gamma}}.
\end{aligned}$$

Then, we have

$$\sum_{\ell=1}^L \sum_{t=0}^{T-1} \Delta_t^\ell(x_t^{*\ell}, s_t^\ell, u_t^{*\ell}) \leq \left(\frac{\delta_\psi T}{\sqrt{\lambda}} + \frac{4\delta_\psi T\sqrt{L}}{\sqrt{\gamma}} \right) \chi(L), \quad (130)$$

Substituting eq. (129) and eq. (130) in eq. (128), we get

$$\mathcal{R}(L) \leq \sigma \sqrt{2LT \log(1/\delta)} + \left(\frac{\delta_\psi T}{\sqrt{\lambda}} + \frac{4\delta_\psi T\sqrt{L}}{\sqrt{\gamma}} \right) \chi(L). \quad (131)$$

the proof follows by substituting $\chi(L)$ defined in eq. (125) in eq. (131), and the probability follows from the union bound.

³Since the system and the weight matrices are known, this assumption can be satisfied by an appropriate choice of the weight matrices.

F Least-squares value iteration for the formulation with process noise in Section 4

In this Appendix, we formulate the regularized least squares regression for the formulation presented in Section 4. Following similar steps as in Appendix B, we define the Bellman target at time t as

$$g_t(x_t, s_t, u_t) = c(x_t, s_t, u_t) + \min_v \widehat{Q}_{t+1}(x_{t+1}, s_{t+1}, v), \quad (132)$$

where x_{t+1} and s_{t+1} denote the states resulting from taking action u_t in states x_t and s_t , and $\widehat{Q}_{t+1}(x_{t+1}, s_{t+1}, v)$ is the estimate of the state-action value function at time $t+1$. We re-write eq. (132) as

$$\begin{aligned} g_t(x_t, s_t, u_t) &= c(x_t, s_t, u_t) + \widehat{V}_{t+1}^*(x_{t+1}, s_{t+1}) \\ &\stackrel{(b)}{=} c(x_t, s_t, u_t) + x_{t+1}^\top G_{t+1} x_{t+1} + 2\widehat{h}_{t+1}^\top(s_{t+1})x_{t+1} + \widehat{q}_{t+1}(s_{t+1}) + \sum_{i=t+2}^T \text{tr}[G_i \Sigma_w], \end{aligned} \quad (133)$$

where in step (b), we have used Theorem 4.1, and $\widehat{h}_{t+1}^\ell(\cdot)$ and $\widehat{q}_{t+1}^\ell(\cdot)$ are as in eq. (49) and eq. (50), respectively. The only modification induced by process noise is the appearance of the constant term $\text{tr}(G_{t+1}\Sigma_w)$. The temporal difference (TD) error is written as

$$\begin{aligned} \varepsilon_t(x_t, s_t, u_t, x_{t+1}, s_{t+1}) &= g_t(x_t, s_t, u_t, x_{t+1}, s_{t+1}) - \widehat{Q}_t(x_t, s_t, u_t) \\ &= c(x_t, s_t, u_t) + \min_v \widehat{Q}_{t+1}(x_{t+1}, s_{t+1}, v) - \widehat{Q}_t(x_t, s_t, u_t) \\ &= c(x_t, s_t, u_t) + x_{t+1}^\top G_{t+1} x_{t+1} + 2\widehat{h}_{t+1}^\top(s_{t+1})x_{t+1} + \widehat{q}_{t+1}(s_{t+1}) + \sum_{i=t+2}^T \text{tr}[G_i \Sigma_w] \\ &\quad - c(x_t, s_t, u_t) - (Ax_t + Bu_t)^\top G_{t+1} (Ax_t + Bu_t) - \phi(s_t)^\top Y(x_t, u_t) \widehat{\theta}_{t+1} - \sum_{i=t+1}^T \text{tr}[G_i \Sigma_w] \\ &= x_{t+1}^\top G_{t+1} x_{t+1} + 2\widehat{h}_{t+1}^\top(s_{t+1})x_{t+1} + \widehat{q}_{t+1}(s_{t+1}) - \phi(s_t)^\top Y(x_t, u_t) \widehat{\theta}_{t+1} \\ &\quad - (Ax_t + Bu_t)^\top G_{t+1} (Ax_t + Bu_t) - \text{tr}(G_{t+1} \Sigma_w), \end{aligned} \quad (134)$$

The TD error $\varepsilon_t(x_t, s_t, u_t, x_{t+1}, s_{t+1})$ in eq. (134) captures the discrepancy between the Bellman target and the current estimate of the Q -function. In Least-Squares Value Iteration (LSVI) in Algorithm 1, we minimize the squared TD error over the dataset collected up to episode $\ell-1$, to obtain an updated estimate of the Q -function. Specifically, the parameters $\widehat{\theta}_t$ at episode ℓ denoted by θ_t^ℓ is obtained by solving the following regularized least-squares problem

$$\begin{aligned} \theta_{t+1}^\ell &= \arg \min_{\theta} \underbrace{\sum_{j=1}^{\ell-1} \varepsilon_t(x_t^j, s_t^j, u_t^j, x_{t+1}^j, s_{t+1}^j)^2}_{J} + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \sum_{j=1}^{\ell-1} \left((Ax_t^j + Bu_t^j)^\top G_{t+1} (Ax_t^j + Bu_t^j) + \phi(s_t^j)^\top Y(x_t^j, u_t^j) \theta + \text{tr}(G_{t+1} \Sigma_w) \right. \\ &\quad \left. - (x_{t+1}^j)^\top G_{t+1} x_{t+1}^j - 2 \left(h_{t+1}^\ell(s_{t+1}^j) \right)^\top x_{t+1}^j - q_{t+1}^\ell(s_{t+1}^j) \right)^2 + \lambda \|\theta\|_2^2. \end{aligned} \quad (135)$$

Taking the derivative of eq. (135) with respect to θ , we get

$$\begin{aligned} \frac{\partial J}{\partial \theta} &= 2 \sum_{j=1}^{\ell-1} Y^\top(x^j, u^j) \phi(s^j) \left((Ax_t^j + Bu_t^j)^\top G_{t+1} (Ax_t^j + Bu_t^j) + \phi(s_t^j)^\top Y(x_t^j, u_t^j) \theta + \text{tr}(G_{t+1} \Sigma_w) \right. \\ &\quad \left. - (x_{t+1}^j)^\top G_{t+1} x_{t+1}^j - 2 \left(h_{t+1}^\ell(s_{t+1}^j) \right)^\top x_{t+1}^j - q_{t+1}^\ell(s_{t+1}^j) \right) + 2\lambda \theta. \end{aligned}$$

Setting the above derivative to zero and solving for θ , we get

$$\begin{aligned}
\theta_{t+1}^\ell &= \Lambda_t^{-1} \sum_{j=1}^{\ell-1} Y^\top(x^j, u^j) \phi(s^j) \epsilon^\ell(x_t^i, x_{t+1}^i, s_{t+1}^i, u_t^i), \\
\Lambda_t &= \sum_{j=1}^{\ell-1} Y(x^j, u^j)^\top \phi(s^j) \phi(s^j)^\top Y(x^j, u^j) + \lambda I_{d(n+1)}, \\
\epsilon^\ell(x_t^i, x_{t+1}^i, s_{t+1}^i, u_t^i) &= (x_{t+1}^i)^\top G_{t+1} x_{t+1}^i - 2(h_{t+1}^\ell(s_{t+1}^i))^\top x_{t+1}^i + q_{t+1}^\ell(s_{t+1}^i) \\
&\quad - (Ax_t^i + Bu_t^i)^\top G_{t+1} (Ax_t^i + Bu_t^i) - \text{tr}(G_{t+1} \Sigma_w).
\end{aligned} \tag{136}$$

G Least-squares value iteration for the setting with unknown dynamics in Section 5

In this appendix, we derive the least-squares value iteration updates for the unknown-dynamics setting in Section 5 and present Algorithm 2, which extends Algorithm 1 to the case where A and B are unknown. Using the expression of Q_t in eq. (24), we can write

$$\begin{aligned}
Q_t^*(x, s, u) &= c(x, s, u) + z^\top P_{t+1} z + 2 \sum_{i=1}^d \phi_i(s) z^\top \tilde{h}_{i,t+1} + \sum_{i=1}^d \phi_i(s) \bar{q}_i(t+1) \\
&= c(x, s, u) + (z^\top \otimes z^\top) \text{vec}(P_{t+1}) + \sum_{i=1}^d \phi_i(s) [2z^\top \quad 1] \begin{bmatrix} \tilde{h}_{i,t+1} \\ \bar{q}_{i,t+1} \end{bmatrix} \\
&= c(x, s, u) + (z^\top \otimes z^\top) \underbrace{\Gamma \text{vech}(P_{t+1})}_{P_{v,t+1}} + \sum_{i=1}^d \phi_i(s) [2z^\top \quad 1] \underbrace{\begin{bmatrix} \tilde{h}_{i,t+1} \\ \bar{q}_{i,t+1} \end{bmatrix}}_{\bar{\theta}_{i,t+1}} \\
&= c(x, s, u) + (z^\top \otimes z^\top) \Gamma P_{v,t+1} + [\phi_1(s) \quad \cdots \quad \phi_d(s)] \underbrace{\begin{bmatrix} [2z^\top \quad 1] & 0 & \cdots & 0 \\ 0 & [2z^\top \quad 1] & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & [2z^\top \quad 1] \end{bmatrix}}_{I_d \otimes [2z^\top \quad 1]} \underbrace{\begin{bmatrix} \bar{\theta}_{1,t+1} \\ \vdots \\ \bar{\theta}_{d,t+1} \end{bmatrix}}_{\bar{\theta}_{t+1}} \\
&= c(x, s, u) + (z^\top \otimes z^\top) \Gamma P_{v,t+1} + \phi(s)^\top (I_d \otimes [2z^\top \quad 1]) \bar{\theta}_{t+1} \\
&= c(x, s, u) + \underbrace{\begin{bmatrix} z^\top \otimes z^\top & \phi(s)^\top \end{bmatrix}}_{\tilde{\Psi}(x,s,u)^\top} \underbrace{\begin{bmatrix} \Gamma & 0 \\ 0 & I_d \otimes [2z^\top \quad 1] \end{bmatrix}}_{\theta_{t+1}} \begin{bmatrix} P_{v,t+1} \\ \bar{\theta}_{t+1} \end{bmatrix} \\
&= c(x, s, u) + \tilde{\Psi}(x, s, u)^\top \theta_{t+1}.
\end{aligned} \tag{137}$$

Following similar steps as in Appendix B, we define the Bellman target at time t as

$$g_t = c(x_t, s_t, u_t) + \underbrace{\min_v \widehat{Q}_{t+1}(x_{t+1}, s_{t+1}, v)}_{y_t}, \tag{138}$$

where x_{t+1} and s_{t+1} denote the states resulting from taking action u_t in states x_{t+1} and s_{t+1} , and $\widehat{Q}_{t+1}(x_{t+1}, s_{t+1}, v)$ is the estimate of the state-action value function at time $t+1$. At time t , the weight parameter θ_{t+1} is obtained by solving the regularized least-squares problem

$$\hat{\theta}_{t+1} = \arg \min_{\theta} \sum_{i=1}^{\ell-1} \left(\tilde{\Psi}(x_t^i, s_t^i, u_t^i)^\top \theta - y_t^i \right)^2 + \lambda \|\theta\|_2^2.$$

Following similar steps as in Appendix B, this problem admits the closed-form solution

$$\Lambda_t^\ell = \sum_{i=1}^{\ell-1} \tilde{\Psi}(x_t^i, s_t^i, u_t^i) \tilde{\Psi}(x_t^i, s_t^i, u_t^i)^\top + \lambda I_{n_\theta}, \quad \text{with } n_\theta = \frac{(n+m+2d)(n+m+1)}{2}, \quad (139)$$

$$\hat{\theta}_{t+1}^\ell = (\Lambda_t^\ell)^{-1} \sum_{i=1}^{\ell-1} \tilde{\Psi}(x_t^i, s_t^i, u_t^i) y_t^i. \quad (140)$$

$$y_t^i = x_{t+1}^i{}^\top G_{t+1}^\ell x_{t+1}^i + 2(h_{t+1}^\ell(s_{t+1}^i))^\top x_{t+1}^i + q_{t+1}^\ell(s_{t+1}^i), \quad (141)$$

where

$$G_{t+1}^\ell = P_{11,t+2}^\ell + W - (P_{12,t+2}^\ell + D)(R + P_{22,t+2}^\ell)^{-1}(P_{12,t+2}^\ell + D)^\top, \quad (142)$$

$$h_{t+1}^\ell(s_{t+1}^i) = \begin{bmatrix} I & K_{x,t+1}^\ell{}^\top \end{bmatrix} (\phi(s_{t+1}^i)^\top \otimes I_{n+m}) \tilde{h}_{t+2}^\ell + (F + K_{x,t+1}^\ell{}^\top H^\top) s_{t+1}^i, \quad (143)$$

$$\begin{aligned} q_{t+1}^\ell(s_{t+1}^i) &= \phi(s_{t+1}^i)^\top \tilde{q}_{t+2}^\ell + s_{t+1}^i{}^\top (M + H K_{s,t+1}^\ell) s_{t+1}^i \\ &\quad - 2s_{t+1}^i{}^\top H (R + P_{22,t+2}^\ell)^{-1} (\phi(s_{t+1}^i)^\top \otimes [0_{m \times n} \quad I_m]) \tilde{h}_{t+2}^\ell \\ &\quad - (\tilde{h}_{t+2}^\ell)^\top \left(\phi(s_{t+1}^i) \otimes \begin{bmatrix} 0_{n \times m} \\ I_m \end{bmatrix} \right) (R + P_{22,t+2}^\ell)^{-1} (\phi(s_{t+1}^i)^\top \otimes [0_{m \times n} \quad I_m]) \tilde{h}_{t+2}^\ell, \end{aligned} \quad (144)$$

$$K_{x,t+1}^\ell = -(R + P_{22,t+2}^\ell)^{-1} (P_{12,t+2}^\ell + D)^\top, \quad (145)$$

$$K_{s,t+1}^\ell = -(R + P_{22,t+2}^\ell)^{-1} H^\top. \quad (146)$$

Note that the half vectorization ensures that P_{t+1} in step 17 of Algorithm 2 is symmetric. However, to enforce the constraint $P_{t+1} \succeq 0$, we project the estimate onto the positive semidefinite cone. Let $P_{t+1} = U \Sigma U^\top$ be the eigenvalue decomposition of the updated matrix. We then replace Σ by $\Sigma_+ = \text{diag}(\max\{\sigma_1, 0\}, \dots, \max\{\sigma_{n+m}, 0\})$, and set $P_{t+1} = U \Sigma_+ U^\top$, as shown in steps 18-20, where $\{\sigma_1, \dots, \sigma_{n+m}\}$ are the eigenvalues of P_{t+1} . In step 22, we apply the control policy in eq. (25) with an added exploration term. Specifically, the input is chosen as

$$u_t^\ell = \hat{K}_{x,t} x_t^\ell + \hat{K}_{s,t} s_t^\ell + \hat{K}_{h,t} (\phi(s_{t+1}^i)^\top \otimes [0_{m \times n} \quad I_m]) \hat{h}_{t+1}^\ell + \eta_t^\ell, \quad (147)$$

where η_t^ℓ is an exploration signal used to ensure sufficient excitation for learning the unknown system parameters.

H Half vectorization

For a matrix $A \in \mathbb{R}^{n \times m}$, let $\text{vec}(A) \in \mathbb{R}^{nm}$ denotes the vectorization of A , which is obtained by stacking the columns of A on top of one another. For a symmetric matrix $A \in \mathbb{R}^{n \times n}$,

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{12} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{nn} \end{bmatrix},$$

we use $\text{vech}(A) \in \mathbb{R}^{n(n+1)/2}$ to denote the half-vectorization of A , which is obtained by vectorizing the lower triangular part of A , i.e.,

$$\text{vech}(A) = [a_{11}, a_{12}, \dots, a_{1n}, a_{22}, \dots, a_{2n}, \dots, a_{nn}]^\top.$$

For a symmetric matrix $A \in \mathbb{R}^{n \times n}$, let $P \in \mathbb{R}^{n^2 \times n(n+1)/2}$ be the unique matrix such that $\text{vec}(A) = P \text{vech}(A)$, where

$$P = \sum_{i \geq j}^n \text{vec}(Q_{ij}) u_{ij}^\top, \quad (148)$$

where $Q_{ij} \in \mathbb{R}^{n \times n}$ is a matrix with 1 in the (i, j) and (j, i) positions and 0 elsewhere, and $u_{ij} \in \mathbb{R}^{n(n+1)/2}$ is a unit vector with 1 in the position $(j-1)n + i - \frac{1}{2}j(j-1)$ and 0 elsewhere.

Algorithm 2 Least-Squares Value Iteration with Unknown System Dynamics

-
- 1: **Given:** episodes L , horizon T , regularizer $\lambda > 0$, projection radius R_θ
 - 2: Initialize $\hat{\theta}_{t+1}^1 \leftarrow 0$ and $\Lambda_t^1 \leftarrow \lambda I_{n_\theta}$, with $n_\theta = \frac{(n+m+2d)(n+m+1)}{2}$, for all $t = 0, \dots, T-1$ (for $\ell = 1$)
 - 3: **for** episode $\ell = 1, \dots, L$ **do**
 - 4: Sample $x_0^\ell \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_x)$ and $s_0^\ell \stackrel{\text{i.i.d.}}{\sim} \mu_0$
 - 5: **Backward pass: parameter updates using data from episodes $1, \dots, \ell-1$**
 - 6: **for** $t = T-1, \dots, 0$ **do**
 - 7: Define dataset $\mathcal{D}^{\ell-1} := \{(x_t^i, s_t^i, u_t^i, x_{t+1}^i, s_{t+1}^i) : i < \ell\}$ (cf. eq. (11))
 - 8: For each sample in $\mathcal{D}^{\ell-1}$, compute the Bellman target

$$y_t^i \leftarrow \min_v \hat{Q}_{t+1}(x_{t+1}^i, s_{t+1}^i, v; \hat{\theta}_{t+2}^\ell)$$

- 9: Form features $\tilde{\Psi}_t^i \leftarrow \tilde{\Psi}(x_t^i, s_t^i, u_t^i)$ (from eq. (137))
- 10: Update Gram matrix and least-squares solution:

$$\Lambda_t^\ell \leftarrow \sum_{i=1}^{\ell-1} \tilde{\Psi}_t^i (\tilde{\Psi}_t^i)^\top + \lambda I_{n_\theta}, \quad \text{where } n_\theta = \frac{(n+m+2d)(n+m+1)}{2},$$

$$\hat{\theta}_{t+1}^\ell \leftarrow (\Lambda_t^\ell)^{-1} \sum_{i=1}^{\ell-1} \tilde{\Psi}_t^i y_t^i,$$

- 11: **if** $\|\hat{\theta}_{t+1}^\ell\|_2 > R_\theta$ **then**
- 12: $\hat{\theta}_{t+1}^\ell \leftarrow \frac{R_\theta}{\|\hat{\theta}_{t+1}^\ell\|_2} \hat{\theta}_{t+1}^\ell$
- 13: **end if**
- 14: **end for**
- 15: **Forward pass: greedy roll-out and data collection**
- 16: **for** $t = 0, \dots, T-1$ **do**
- 17: Extract \hat{P}_{t+1}^ℓ and \hat{h}_{t+1}^ℓ from $\hat{\theta}_{t+1}^\ell$
- 18: $P_{t+1} = U \Sigma U^\top$
- 19: $\Sigma_+ = \text{diag}(\max\{\sigma_1, 0\}, \dots, \max\{\sigma_{n+m}, 0\})$
- 20: $P_{t+1} \leftarrow U \Sigma_+ U^\top$
- 21: Form gains using eq. (26) with blocks $\hat{P}_{21,t+1}^\ell, \hat{P}_{22,t+1}^\ell$:

$$\hat{K}_{x,t} = -(R + \hat{P}_{22,t+1}^\ell)^{-1} (D^\top + \hat{P}_{21,t+1}^\ell), \quad \hat{K}_{s,t} = -(R + \hat{P}_{22,t+1}^\ell)^{-1} H^\top, \quad \hat{K}_{h,t} = -(R + \hat{P}_{22,t+1}^\ell)^{-1}$$

- 22: Apply optimal policy with exploration term η_t^ℓ (cf. eq. (25)):

$$u_t^\ell \leftarrow \hat{K}_{x,t} x_t^\ell + \hat{K}_{s,t} s_t^\ell + \hat{K}_{h,t} (\phi(s_{t+1}^i)^\top \otimes [0_{m \times n} \quad I_m]) \hat{h}_{t+1}^\ell + \eta_t^\ell$$

- 23: Take action u_t^ℓ , observe x_{t+1}^ℓ and s_{t+1}^ℓ
 - 24: **end for**
 - 25: **end for**
-