
Taming “data-hungry” reinforcement learning? Stability in continuous state-action spaces

Yaqi Duan

Department of Technology, Operations, and Statistics
Stern School of Business, New York University
New York, NY 10012
yaqi.duan@stern.nyu.edu

Martin J. Wainwright

Laboratory for Information and Decision Systems, Statistics and Data Science Center
Department of Electrical Engineering and Computer Science, and Department of Mathematics
Massachusetts Institute of Technology Cambridge, MA 02139
wainwrigwork@gmail.com

Abstract

We introduce a novel framework for analyzing reinforcement learning (RL) in continuous state-action spaces, and use it to prove fast rates of convergence in both off-line and on-line settings. Our analysis highlights two key stability properties, relating to how changes in value functions and/or policies affect the Bellman operator and occupation measures. We argue that these properties are satisfied in many continuous state-action Markov decision processes. Our analysis also offers fresh perspectives on the roles of pessimism and optimism in off-line and on-line RL.

1 Introduction

Many domains of science and engineering involve making a sequence of decisions over time, with previous decisions influencing the future in uncertain ways [1, 13, 22, 31, 32]. For instance, clinicians managing diabetes [36] or engineers optimizing plasma control in tokamak systems [5] must develop policies that adapt based on evolving conditions and lead to desirable outcomes over a longer period. Markov decision processes (MDPs) and reinforcement learning (RL) provide frameworks and methods for estimating effective policies for such sequential problems. While RL excels in data-rich scenarios such as competitive gaming (e.g., AlphaGo and its extensions [30]), its application in data-scarce areas like healthcare [36] and finance [27] remains challenging due to lack of history, or underlying non-stationarity. With limited data, characterizing and improving the *sample complexity* of RL methods becomes critical.

Considerable research effort has been devoted to studying RL sample complexity in many settings. Existing studies for either the generative or the off-line settings (e.g., [21, 37, 35]) give procedures that, when applied to a dataset of size n , yield a value gap that decays at the rate $1/\sqrt{n}$. In the on-line setting, there are various procedures that yield cumulative regret that grows at the rate \sqrt{T} (e.g., [18, 20, 19, 6]). In contrast, the main result of this paper is to formalize conditions, suitable for RL in continuous domains, under which *much faster rates can be obtained using the same dataset*, achieving a value gap decay of $1/n$ and reducing regret growth to $\log T$.

As revealed by our analysis, these accelerated rates depend on certain *stability properties*, ones that—as we argue—are naturally satisfied in many control problems with continuous state-action

spaces. Roughly speaking, these conditions ensure that the evolution of the dynamic system depends in a “smooth” way on the influence of decision policy. Such notions of stability should be expected in various controlled systems with continuous state-action spaces. In robotics, for example, a minor torque or motion perturbation that occurs during a single step should not cause a notable deviation from the intended trajectory. Similarly, in clinical treatment, slight deviations in medication dosage should not significantly compromise effectiveness or safety.

1.1 A simple illustrative example: Mountain Car

The “Mountain Car” problem, a benchmark continuous control task, illustrates the acceleration phenomenon and underlying stability. In this task, as shown in Figure 1(a), a car must reach the top of a hill by adjusting its acceleration within the interval $[-1, 1]$. We employed fitted Q -iteration (FQI) with carefully selected linear basis functions to derive near-optimal policies with off-line data. This learning procedure exhibits a value sub-optimality decay at a rate of $1/n$, a significant improvement over the classical rate of $1/\sqrt{n}$, as detailed in Figure 1(b). (See Appendix D for further explanation. The experiment ran for 3 days on two laptops, each equipped with an Apple M2 Pro CPU and 16 GB RAM.) In this example, slight perturbations in the driving policy lead to only modest changes in future trajectories, which shows the stability. Our theoretical analysis confirms that fast rates are achievable in this and similar continuous control tasks when such stability properties are present.

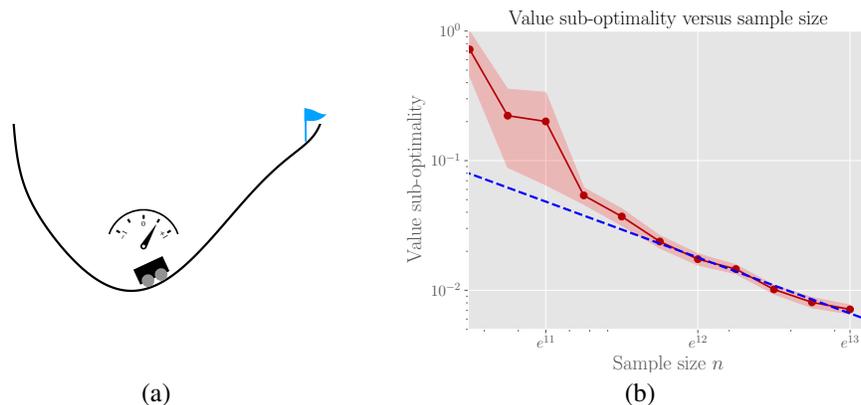


Figure 1: Illustration of the “fast rate” phenomenon using FQI on the Mountain Car problem. Each **red point** in the plot represents the average value sub-optimality $J(\pi^*) - J(\hat{\pi}_n)$ from $T = 80$ Monte Carlo trials, with the shaded area showing twice the standard errors. The **blue dashed line** is a least-squares fit to the last 6 data points, yielding a 95% confidence interval of $(-1.084, -0.905)$ for the slope, significantly faster than the typical -0.5 “slow rate”.

1.2 Contributions of this paper

With this high-level perspective in mind, let us summarize the key contributions of this paper.

Fast rate of convergence: We develop a framework for analyzing RL in continuous state-action spaces, and use it to prove a general result (Theorem 1) under which fast rates can be obtained. The key insight is that stability conditions lead to upper bounds on the value sub-optimality that are proportional to the *squared* norm of Bellman residuals. In the off-line setting, this quadratic scaling improves convergence from a rate of $n^{-\frac{1}{2}}$ to n^{-1} , while in on-line learning, it enhances the regret bound from \sqrt{T} to $\log T$.

Reconsidering pessimism and optimism principles: Our framework provides a novel perspective on the roles of pessimism [21, 4] and optimism [18, 20, 19, 6, 12] in off-line and on-line RL. Our theory reveals that there are settings in which *neither pessimism nor optimism* are required for effective policy optimization—in particular, they are not required as long as one has a sufficiently accurate pilot estimate policy. Moreover, our analysis shows that some procedures based on certainty

equivalence can achieve fast-rate convergence, showing that the benefits gained from incorporating additional pessimism or optimism measures may be limited in this context.

1.3 Related work

In this section, we discuss related work having to do with fast rates in optimization and statistics.

Fast rates in stochastic optimization and risk minimization: For many statistical estimators (e.g., likelihood methods, empirical risk minimization), it is well-understood that the local geometry around the optimum determines whether fast rates can be obtained. For instance, when the loss function exhibits some form of strong convexity (such as exp-concave loss) or strict saddle properties, it can lead to significant reductions in additive regret from $\mathcal{O}(\sqrt{T})$ to just $\mathcal{O}(\log T)$ in stochastic approximation (e.g., [15]), or a decrease in the error rate from $n^{-\frac{1}{2}}$ to n^{-1} in empirical risk minimization [23, 14]. These fast rate phenomena rely on a form of stability, one which relates the similarity of functions to the closeness of their optima. Our work develops a new framework for analyzing value-based RL methods, focusing on identifying specific stability conditions and inherent curvature properties that promote fast rate convergence in RL, similar to the role of stability analysis in statistical learning.

Fast rates in reinforcement learning: In the RL literature, there are various lines of work related to fast rates, but the underlying mechanisms are typically different from those considered here. For problems with discrete state-action spaces, there is a line of recent work [17, 16, 33, 25] that performs gap/marginal-dependent analyses of RL algorithms. However, such separation assumptions are not helpful for continuous action spaces. Other work for discrete state-action spaces [28] has shown convergence rates in off-line RL are influenced by data quality, with a nearly-expert dataset enabling faster rate. In contrast, our analysis reveals that for off-line RL in continuous domains, fast convergence can occur whether or not the dataset has good coverage properties.

An important sub-class of continuous state-action problems are those with linear dynamics and quadratic reward functions (LQR for short). For such problems, it has been shown [24, 29] that value sub-optimality can be connected with the squared error in system identification. Our general theory can also be used to derive guarantees for LQR problems, as we explore in more detail in a follow-up paper [8]. Stability also arises in the analysis of (deterministic) policy optimization and Newton-type algorithms [26, 3], where it is possible to show superlinear convergence in a local neighborhood. This accelerated rate stems from the smoothness of the on-policy transition operator \mathcal{P}^{π_f} with respect to changes in the value function f ; for instance, see condition (10) in Puterman and Brumelle [26]. Our framework exploits related notions of smoothness, but is tailored to the stochastic setting of reinforcement learning, in which understanding the effect of function approximation and finite sample sizes is essential.

2 Fast rates for value-based reinforcement learning

Let us now set up and state the main result of this paper. We begin in Section 2.1 with background on Markov decision processes (MDPs) and value-based methods, before turning to the statement of our main result in Section 2.2. In Section 2.3, we provide intuition for why stability leads to faster rates, and discuss consequences for both the off-line and on-line settings of RL.

2.1 Markov decision processes and value-based methods

Basic set-up: We consider an episodic Markov decision process (MDP) defined by a quadruple $(\mathcal{S}, \mathcal{A}, \mathcal{P} = \{\mathcal{P}_h\}_{h=1}^{H-1}, \{r_h\}_{h=1}^H)$. We assume that the rewards $r_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ are known; however, this condition can be relaxed. A policy π_h at time h is a mapping from any state s to a distribution $\pi_h(\cdot | s)$ over the action space \mathcal{A} . If the support of $\pi_h(\cdot | s)$ is a singleton, we also let $\pi_h(s) \in \mathcal{A}$ denote the single action to be chosen at state s . Given an initial distribution ξ_1 over the states at time $h = 1$, the *expected reward* obtained by choosing actions according to a policy sequence $\pi = (\pi_1, \dots, \pi_H)$ is given by $J(\pi) \equiv J(\pi; \xi_1) := \mathbb{E}_{\xi_1, \pi} [\sum_{h=1}^H r_h(S_h, A_h)]$, where $S_1 \sim \xi_1$, $S_{h+1} \sim \mathcal{P}_h(\cdot | S_h, A_h)$ and $A_h \sim \pi_h(\cdot | S_h)$ for $h = 1, 2, \dots, H$. Our goal is to estimate an *optimal policy* $\pi^* \in \arg \max_{\pi} J(\pi)$.

Value functions and Bellman operators: Starting from a given state-action pair (s, a) at stage h , the expected return over subsequent stages defines the *state-action value function* $Q_h^\pi(s, a) := \mathbb{E}_\pi [\sum_{h'=h}^H r(S_{h'}, A_{h'}) \mid S_h = s, A_h = a]$. The sequence of functions $\mathbf{Q}^\pi = (Q_1^\pi, \dots, Q_H^\pi)$ known as the *Q-functions* associated with π .

The *Q-functions* \mathbf{Q}^π have an important connection with the *Bellman evaluation operator* for π . For any policy π and stage h , we introduce a linear transition operator $(\mathcal{P}_h^\pi f)(s, a) := \int_{\mathcal{S} \times \mathcal{A}} f(s', a') \mathcal{P}_h(ds' \mid s, a) \pi_{h+1}(da' \mid s')$ for any function $f \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. With this notation, the *Bellman evaluation operator* at stage h takes the form

$$(\mathcal{T}_h^\pi f)(s, a) := r_h(s, a) + (\mathcal{P}_h^\pi f)(s, a). \quad (1)$$

From classical dynamic programming, the *Q-functions* \mathbf{Q}^π must satisfy the Bellman relations $Q_h^\pi(s, a) = (\mathcal{T}_h^\pi Q_{h+1}^\pi)(s, a)$ for $h = 1, \dots, H - 1$.

Bellman principle for optimal policies: Under mild regularity conditions, there is at least one policy π^* such that, for any other policy π , we have $Q_h^{\pi^*}(s, a) \geq Q_h^\pi(s, a)$, for any $h \in [H]$, and uniformly over all state-action pairs (s, a) . Any optimal policy π^* must be greedy with respect to the optimal *Q-function* \mathbf{Q}^* . By classical dynamic programming, the optimal *Q-function* \mathbf{Q}^* is obtained by setting $Q_H^* = r_H$, and then recursively computing $Q_h^* = \mathcal{T}_h^* Q_{h+1}^*$ for $h = H - 1, \dots, 2, 1$, with the *Bellman optimality operator* defined as

$$(\mathcal{T}_h^* f)(s, a) := r_h(s, a) + \mathbb{E}_h \left[\max_{a' \in \mathcal{A}} f(S', a') \mid s, a \right] \quad \text{for } S' \sim \mathcal{P}_h(\cdot \mid s, a). \quad (2)$$

Value-based RL methods The main result of this paper applies to a broad class of methods for reinforcement learning. They are known as *value-based*, due to their reliance on the following two step approach for approximating an optimal policy π^* : (1) Construct an estimate $\hat{\mathbf{Q}} = (\hat{f}_1, \dots, \hat{f}_H)$ of the optimal value function $\mathbf{Q}^* = (Q_1^*, \dots, Q_H^*)$. (2) Use $\hat{\mathbf{Q}}$ to compute the greedy-optimal policy $\hat{\pi}_h(s) \in \arg \max_a \hat{f}_h(s, a)$ for $h = 1, 2, \dots, H$. It should be noted that there is considerable freedom in the design of a value-based method, since different methods can be used to approximate value functions in Step 1. Rather than applying to a single method, our main result applies to a very broad class of these methods.

Underlying any value-based method is a class \mathcal{F} of functions $(s, a) \mapsto f(s, a)$ used to approximate the state-action value functions.¹ We assume that the function class \mathcal{F} is rich enough—relative to the Bellman evaluation operators—to ensure that for any greedy policy π induced by some $\mathbf{f} = (f_1, \dots, f_H) \in \mathcal{F}^H$, we have the inclusion $\mathcal{T}_h^\pi \mathcal{F} \subseteq \mathcal{F}$ for $h = 1, \dots, H - 1$. We see that this condition depends on the structure of the transition distributions $\mathcal{P}_h(\cdot \mid s, a)$. In many practical examples, the reward function itself has some number of derivatives, and these transition distributions perform some type of smoothing, so that we expect that the output of the Bellman update, given a suitably differentiable function, will remain suitably differentiable.

2.2 Stable problems have fast rates

We now turn the central question in understanding the behavior of any value-based method:

*How to translate “closeness” of the Q-function estimate $\hat{\mathbf{Q}}$
to a bound on the value gap $J(\pi^*) - J(\hat{\pi})$?*

At a high level, existing theory provides guarantees of the following type: if the *Q-function* estimates are ε -accurate for some $\varepsilon \in (0, 1)$, then the value gap is bounded by a quantity proportional to ε . In contrast, our main result shows that when the MDP is stable in a suitable sense, the value gap can be upper bounded by a quantity proportional to ε^2 . This *quadratic as opposed to linear scaling* encapsulates the “fast rate” phenomenon of this paper.

Our analysis isolates two key stability properties required for faster rates; both are Lipschitz conditions with respect to a certain norm. Here we define them with respect to the L^2 -norm induced by the

¹In general, different function classes may be selected at each stage $h = 1, 2, \dots, H$; here, so as to reduce notational clutter, we assume that the same function class \mathcal{F} is used for each stage.

state-action occupation measure induced by the optimal policy—namely

$$\|f\|_h := \sqrt{\mathbb{E}_{\pi^*}[f^2(S_h, A_h)]} \quad \text{for any } f \in \partial\mathcal{F}^2, \quad (3)$$

and over a neighborhood \mathcal{N} of the optimal Q -value function Q^* .

Bellman stability: The first condition measures the stability of the Bellman optimality operator (2): in particular, we require that there is a scalar κ_h^* such that

$$\|\mathcal{T}_h^* f_{h+1} - \mathcal{T}_h^* Q_{h+1}^*\|_h \leq \kappa_h^* \|f_{h+1} - Q_{h+1}^*\|_{h+1} \quad (\text{Stb}(\mathcal{T}))$$

for any $f \in \mathcal{N}$. Moreover, for any pair (h, h') of indices such that $1 \leq h < h' \leq H - 1$, we define

$$\kappa_{h,h'}(\mathcal{T}^*) := \kappa_h^* \kappa_{h+1}^* \cdots \kappa_{h'-1}^*.$$

Condition $(\text{Stb}(\mathcal{T}))$ is directly linked to the stability of estimating the Q -function Q^* . In typical estimation procedures, such as approximate dynamic programming, the estimation is carried out iteratively in a backward manner, so that it is important to control the propagation of estimation errors across the iterations. Condition $(\text{Stb}(\mathcal{T}))$ captures this property, since it implies that

$$\|\mathcal{T}_h^* \mathcal{T}_{h+1}^* \cdots \mathcal{T}_{h'-1}^* f_{h'} - \mathcal{T}_h^* \mathcal{T}_{h+1}^* \cdots \mathcal{T}_{h'-1}^* Q_{h'}^*\|_h \leq \kappa_{h,h'}(\mathcal{T}^*) \cdot \|f_{h'} - Q_{h'}^*\|_{h'},$$

which shows how the estimation error $(f_{h'} - Q_{h'}^*)$ at step h' can be controlled in terms of estimation error at an earlier time step $h \leq h'$.

Occupation measure stability: Our second condition is more subtle, and is key in our argument. Let us begin with some intuition. Consider two sequences of policies

$$(\pi_1^*, \dots, \pi_{h-1}^*, \pi_h^*, \pi_{h+1}^*, \dots, \pi_{h'}^*) \quad \text{and} \quad (\pi_1^*, \dots, \pi_{h-1}^*, \pi_h, \pi_{h+1}^*, \dots, \pi_{h'}^*)$$

that only differ at the h -th step, where π_h^* has been replaced by π_h . These two policy sequences induce Markov chains whose distributions differ from stage h onwards, and our second condition controls this difference in terms of the difference $\|f_h - Q_h^*\|_h$ between the two Q -functions f_h and Q_h^* that induce π_h and π_h^* , respectively.

We adopt \mathcal{P}_h^* as a convenient shorthand for the transition operator $\mathcal{P}_h^{\pi^*}$, and define the multi-step transition operator $\mathcal{P}_{h,h'}^* := \mathcal{P}_h^* \mathcal{P}_{h+1}^* \cdots \mathcal{P}_{h'-1}^*$. Using this notation, for any $h' \geq h + 1$, we require that there is a scalar $\kappa_{h,h'}(\pi^*)$ such that

$$\sup_{\substack{g \in \partial\mathcal{F} \\ \|g\|_{h'} > 0}} \frac{|\mathbb{E}_{\pi^*}[(\mathcal{P}_{h,h'}^* g)(S_h, \pi_h^*(S_h)) - (\mathcal{P}_{h,h'}^* g)(S_h, \pi_h(S_h))]|}{\|g\|_{h'}} \leq \kappa_{h,h'}(\pi^*) \frac{\|f_h - Q_h^*\|_h}{\|Q_h^*\|_h} \quad (\text{Stb}(\xi))$$

for any $f \in \mathcal{N}$. The renormalization in this definition serves to enforce a natural scale invariance.

With these notions of stability in hand, we are now equipped to state our main result. Taking as input a value function estimate \hat{Q} , it relates the induced value gap to the *Bellman residuals* $\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h$. Note that these residuals are a way of quantifying proximity to the optimal value function Q^* , which has Bellman residual zero by definition. We assume that \hat{Q} has Bellman residuals bounded as

$$\|\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h\|_h \leq \varepsilon_h \quad \text{for } h = 1, 2, \dots, H - 1 \quad (4a)$$

for some sequence $\varepsilon = (\varepsilon_1, \dots, \varepsilon_{H-1}, \varepsilon_H = 0)$ that satisfies the constraint

$$\varepsilon_h \geq \frac{1}{H-h} \sum_{h'=h+1}^H \varepsilon_{h'} \quad \text{for } h = 1, 2, \dots, H - 1. \quad (4b)$$

This last condition means that the Bellman residual ε_h is larger than or equal to the average of the bounds established after step $h + 1$. It is natural because estimating at step h is at least as challenging as a stage $h' > h$; indeed, any such state h' occurs earlier in the dynamic programming backward iteration process. As a special case, the bound (4b) holds when $\varepsilon_h = \varepsilon$ for all stages.

With this set-up, we have the following guarantee in terms of the stability coefficients $\kappa_{h,h'}(\pi^*)$ and $\kappa_{h,h'}(\mathcal{T}^*)$ from conditions $(\text{Stb}(\xi))$ and $(\text{Stb}(\mathcal{T}))$.

²We let $\partial\mathcal{F}$ be the set of all difference functions of the form $g = f - \tilde{f}$ for some $f, \tilde{f} \in \mathcal{F}$.

Theorem 1. *There is a neighborhood of Q^* such that for any value function estimate \hat{f} with ε -bounded Bellman residuals (4a), the induced greedy policy $\hat{\pi}$ has value gap bounded as*

$$J(\pi^*) - J(\hat{\pi}) \leq 2 \sum_{h=1}^{H-1} \frac{1}{\|Q_h^*\|_h} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\pi^*) \varepsilon_{h'} \right\} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\mathcal{T}^*) \varepsilon_{h'} \right\}. \quad (5)$$

See Appendix A for the proof.

Treating dependence on the stability coefficients as constant, the main take-away is that value suboptimality is bounded above by a quantity proportional to the *squared* norm of the Bellman residuals. Concretely, if the Bellman residuals are uniformly upper bounded by some ε , then equation (5) leads to an upper bound of the form

$$J(\pi^*) - J(\hat{\pi}) \leq c H^3 \varepsilon^2,$$

where c is a universal constant. Due to the quadratic scaling in the Bellman residual error ε , this bound is substantially tighter than the linear in ε rates afforded by a conventional analysis.

2.3 Intuition for fast rates: Smoothness and cancelling terms in the telescope bound

Why does “fast rate” phenomenon formalized in Theorem 1 arise? The fast rates proved in this paper are established by a novel argument, starting from a known telescope bound, which we begin by stating. Given a Q -function estimate $\hat{f} = (\hat{f}_1, \dots, \hat{f}_H)$, let $\hat{\pi}$ denote the induced greedy policy. Then the value gap of $\hat{\pi}$ with respect to an arbitrary comparator policy π is bounded as

$$J(\pi) - J(\hat{\pi}) \leq \sum_{h=1}^{H-1} (\mathbb{E}_\pi - \mathbb{E}_{\hat{\pi}}) [(\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h)]. \quad (6)$$

This result follows by a “telescope” relation induced by the structure of the Bellman updates.³ For completeness, we provide a proof of the telescope bound in Appendix E.2.

A key feature of inequality (6) is the difference of two expectations $\mathbb{E}_\pi - \mathbb{E}_{\hat{\pi}}$, corresponding to the occupation measures under π versus $\hat{\pi}$. In standard uses of this inequality, an initial argument is used to guarantee that one of these expectations is negative, and so can be dropped [20, 21].

In contrast, the proof of our Theorem 1 exploits a more refined approach, one that handles the difference of expectations directly. Doing so can be beneficial—and lead to “fast rates”—because various terms in this difference can cancel each other out. Specifically, under the smoothness conditions that underlie Theorem 1, when applying the telescope inequality (6) with comparator $\pi = \pi^*$, we show that the discrepancy between the occupation measures associated with π^* and $\hat{\pi}$ is of the *same order* as the Bellman residual associated with \hat{f} . Note that the Bellman residuals of \hat{f} already appear on the right-hand side of inequality (6), so that this fortuitous cancellation can be exploited—along with a number of auxiliary results laid out in the proof—so as to upper bound the value gap by a quantity proportional to the squared Bellman residual ε^2 .

It is worthwhile making an explicit comparison of our cancellation approach with the more standard uses of the telescope relation, which typically consider only one portion of the Bellman residuals (e.g., [18, 20, 21, 19, 6, 12, 35]). We do so in the following two subsections.

2.3.1 Pessimism for off-line RL

In the off-line instantiation of RL, the goal is to learn a “good” policy based on a pre-collected dataset \mathcal{D} . Note that no further interaction with the environment is permitted, hence the notion of the learning being off-line. More precisely, an *off-line dataset* \mathcal{D} of size n consists of quadruples

$$\mathcal{D} = \left\{ (s_{h,i}, a_{h,i}, s'_{h,i}, r_{h,i}) \right\}_{i=1}^n,$$

where $s_{h,i}$ and $a_{h,i}$ represent the i -th state and action at the h -th step in the MDP; $s'_{h,i}$ is the successive state; and $r_{h,i} = r_h(s_{h,i}, a_{h,i})$ denotes the scalar reward. Note that while the successive

³Results of this type are known; for example, analogous results can be found in past work (e.g., Theorem 2 of the paper [34]; or Lemma 3.2 in the paper [7]).

states are defined by transition dynamics, and the rewards by the reward function, there are no restrictions on how the state-action pairs $(s_{h,i}, a_{h,i})$ are collected. That is, they need not have been generated by any fixed policy, but may have collected from some ensemble of behavioral policies, or even adaptively by human experts. The goal of off-line reinforcement learning is to use the n -sample dataset \mathcal{D} so as to estimate a policy $\hat{\pi} \equiv \hat{\pi}_n$ that (approximately) maximizes the expected return $J(\hat{\pi}_n)$. We expect that—at least for a sensible method for estimating $\hat{\pi}_n$ —the value gap $J(\pi^*) - J(\hat{\pi}_n)$ should decay to zero as n increases to infinity, and we are interested in understanding this rate of decay.

The use of pessimism is standard in off-line RL algorithms. Its purpose is to mitigate risks associated with “poor coverage” of the off-line dataset. For instance, the naive approach of simply maximizing Q -function estimates based on an off-line dataset can behave poorly when certain portions of the state-action space are not well covered by the given dataset. The pessimism principle suggests to form a *conservative estimate* of the value function—say with

$$\hat{f}_h(s, a) \leq \mathcal{T}_h^* \hat{f}_{h+1}(s, a) \quad (7a)$$

with high probability over state-action pairs (s, a) . Thus, the estimated value $\hat{f}_h(s, a)$ is an underestimate of the Bellman update, a form of conservatism that protects against unrealistically high estimates due to poor coverage. Doing so in the appropriate way ensures that

$$-\mathbb{E}_{\hat{\pi}}[(\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h)] \leq 0. \quad (7b)$$

Applying this upper bound to the inequality (6) yields the sub-optimality bound

$$J(\pi) - J(\hat{\pi}) \leq \sum_{h=1}^{H-1} \mathbb{E}_{\pi}[(\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h)].$$

Upper bounds derived in this manner only contain one portion of the Bellman residual. When the value functions are approximated in a parametric way (e.g., tabular problems, linear function approximation), this line of analysis leads to value sub-optimality decaying at a “slow” $1/\sqrt{n}$ rate in terms of the sample size n (e.g., [21]). In contrast, an application of Theorem 1 can lead to value gaps bounded by $1/n$.

2.3.2 Optimism in on-line RL

In the setting of on-line RL, a learning agent interacts with the environment in a sequential manner, receiving feedback in the form of rewards based on its actions. At the beginning, the learner possesses no prior knowledge of the system’s dynamics. In the t -th episode, the agent learns an optimal policy $\hat{\pi}^{(t)}$ using existing observations, implements the policy and collects data $\{(s_h^{(t)}, a_h^{(t)}, r_h^{(t)})\}_{h=1}^H$ from the new episode. In each round, the system starts at an initial state $s_1^{(t)}$ independently drawn from a fixed distribution ξ_1 .

In this on-line setting, it is common to measure the performance of an algorithm by comparing it, over the T rounds of learning, with an oracle that knows and implements an optimal policy. At each round t , we incur the *instantaneous regret* $J(\pi^*) - J(\hat{\pi}^{(t)})$, where π^* is any optimal policy. Over T rounds, we measure performance in terms of the *cumulative regret*

$$\text{Regret}(\{\hat{\pi}^{(t)}\}_{t=1}^T) := \max_{\text{policy } \pi} \sum_{t=1}^T \{J(\pi) - J(\hat{\pi}^{(t)})\} = \sum_{t=1}^T \underbrace{\{J(\pi^*) - J(\hat{\pi}^{(t)})\}}_{\text{Regret at round } t}. \quad (8)$$

In a realistic problem, the cumulative regret of any procedure grows with T , and our goal is to obtain algorithms whose regret grows as slowly as possible.

In contrast to off-line RL, the on-line setting allows for exploring state-action pairs that have been rarely encountered; doing so makes sense since they might be associated with high rewards. Principled exploration of this type can be effected via the *optimism principle*: one constructs function estimates such that

$$\hat{f}_h(s, a) \geq \mathcal{T}_h^* \hat{f}_{h+1}(s, a) \quad (9a)$$

with high probability over state-action pairs.⁴ Note that $\widehat{f}_h(s, a)$ is optimistic in the sense that it is an over-estimate of the Bellman update $\mathcal{T}_h^* \widehat{f}_{h+1}(s, a)$. In this way, we can ensure that

$$\mathbb{E}_\pi [(\mathcal{T}_h^* \widehat{f}_{h+1} - \widehat{f}_h)(S_h, A_h)] \leq 0. \quad (9b)$$

Combining this inequality with the telescope bound (6) allows one to upper bound the regret as

$$\text{Regret}(\{\widehat{\pi}^{(t)}\}_{t=1}^T) = \sum_{t=1}^T \{J(\pi^*) - J(\widehat{\pi}^{(t)})\} \leq \sum_{t=1}^T \sum_{h=1}^{H-1} \mathbb{E}_{\widehat{\pi}^{(t)}} [(\widehat{f}_h - \mathcal{T}_h^* \widehat{f}_{h+1})(S_h, A_h)].$$

which only includes a single portion of the Bellman residual. In the case of tabular or linear representations of the Q -functions, it results in a regret rate of \sqrt{T} (e.g., see the papers [18, 20]). In contrast, an appropriate use of Theorem 1 leads to regret growing only as $\log(T)$, which corresponds to a much better guarantee.

In summary, then, the fast rates obtained in this paper are based on a different approach than the standard pessimism or optimism principles. Since we deal directly with the difference of expectations in the bound (6), there is no need to nullify either of them through the use of these principles. However, it should be noted that we are assuming smoothness conditions that allow us to control this difference. As we discuss in the sequel, such smoothness conditions rule out certain “hard instances” used in past work on lower bounds (e.g. [18, 20, 21, 37]).

3 Consequences for linear function approximation

In this section, we explore some consequences of our general theory when applied to value-based methods using (finite-dimensional) linear function approximation.

Let $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a given feature map on the state-action space, and consider linear expansions of the form $f_w(s, a) = \langle \phi(s, a), \mathbf{w} \rangle \equiv \sum_{j=1}^d w_j \phi_j(s, a)$ where $\mathbf{w} \in \mathbb{R}^d$ is a weight vector. We adopt the conventional assumption that $\|\phi(s, a)\|_2 \leq 1$ and $r_h(s, a) \in [0, 1]$ for all state-action pairs. Defining the linear function class $\mathcal{F} := \{f_w \mid \mathbf{w} \in \mathbb{R}^d\}$, we note that the Minkowski difference class $\partial \mathcal{F}$ is equal to \mathcal{F} .

In our analysis of linear approximation, we make use of the norm $\|f\|_h := \sqrt{\mathbb{E}_{\pi^*}[f^2(S_h, A_h)]}$, corresponding to L^2 -norm under the occupation measure induced by the optimal policy π^* .

3.1 Consequences for off-line RL

We now turn to some implications of Theorem 1 for off-line reinforcement learning. Let us recall the off-line setting: for each $h = 1, \dots, H - 1$, we are given a dataset $\mathcal{D}_h = \{(s_{h,i}, a_{h,i}, s'_{h,i}, r_{h,i})\}_{i=1}^n$ of quadruples, from which we can compute estimates $\widehat{f} = (\widehat{f}_h)_{h=1}^H$ with certain Bellman residuals $\{\varepsilon_h\}_{h=1}^{H-1}$, which then appear in the bound (5). The remaining factors on the right-hand side of inequality (5) do not depend on the dataset itself (but rather on structural properties of the MDP). Consequently, in terms of statistical understanding, the main challenge is to establish high-probability bounds on the Bellman residuals $\{\varepsilon_h\}_{h=1}^{H-1}$ for a particular estimator.

3.1.1 Fitted Q -iteration (FQI)

As an illustration, let us analyze the use of *fitted Q -iteration* (FQI) for computing estimates of the Q -function. At a given stage $h = 1, \dots, H - 1$, we can use the associated data \mathcal{D}_h to define a regularized objective function

$$\mathcal{L}_h(f, g) := \frac{1}{|\mathcal{D}_h|} \left[\sum_{(s_{h,i}, a_{h,i}, s'_{h,i}, r_{h,i}) \in \mathcal{D}_h} \{f(s_{h,i}, a_{h,i}) - (r_{h,i} + \max_{a \in \mathcal{A}} g(s'_{h,i}, a))\}^2 \right] + \Lambda_h^2(f).$$

Here g represents the target function from stage $h + 1$, and it defines the targeted responses $y_{h,i}(g) := r_{h,i} + \max_{a \in \mathcal{A}} g(s'_{h,i}, a)$. For a given target g , we obtain a Q -function estimate for

⁴Please refer to, for example, Lemma B.3 in the paper [20] for further details.

stage h by minimizing the functional $f \mapsto \mathcal{L}_h(f, g)$. Given that our objective is defined with a quadratic cost, doing so can be understood as a regression method for estimating the conditional expectation that underlies the Bellman update—viz. $\mathcal{T}_h^* g(s, a) = \mathbb{E}[y_{h,i}(g) \mid (s_{h,i}, a_{h,i}) = (s, a)]$. Here $\Lambda_h^2(f) = \lambda_h \|\mathbf{w}\|_2^2$ for $f = \langle \phi(\cdot), \mathbf{w} \rangle$ is a regularizer, with $\lambda_h \geq 0$ being the regularization weight. Given this set-up, we can generate a Q -function estimate $\hat{\mathbf{f}} = (\hat{f}_1, \dots, \hat{f}_H)$ by first initializing $\hat{f}_H = r_H$, and then recursively computing $\hat{f}_h = \arg \min_{f \in \mathcal{F}} \mathcal{L}_h(f, \hat{f}_{h+1})$, for $h = H - 1, H - 2, \dots, 2, 1$.

3.1.2 Fast rates for FQI-based estimates

In the analysis here, we assume that the dataset consists of i.i.d. tuples (but this can be relaxed as needed). We now state a corollary of Theorem 1, applicable to value function estimates based on FQI with ridge regression.

Corollary 1 (Fast rates for ridge-based FQI). *For FQI based on ridge regression, with a sufficiently large sample size n and with suitable choices of the regularization parameters $\{\lambda_h\}_{h=1}^{H-1}$, the bound (5) from Theorem 1 holds with*

$$\varepsilon_h = c \sqrt{\{d(H-h)/n\} \log(dH/\delta)} \quad (10)$$

with probability at least $1 - \delta$.

We omit the proof of Corollary 1, as it follows from standard ridge regression analysis.

Fast rates and comparisons to past work: So as to be able to compare with results from past work, let us consider some consequences of the bound (10) under the following assumptions: (i) $\kappa_{h,h'}(\mathcal{T}^*) = \mathcal{O}(1)$; (ii) $\kappa_{h,h'}(\boldsymbol{\pi}^*) = \mathcal{O}(\sqrt{d})$; (iii) $\|Q_h^*\|_h \asymp H - h + 1$. Then it can be shown that the bound from Corollary 1 takes the form

$$J(\boldsymbol{\pi}^*) - J(\hat{\boldsymbol{\pi}}) \leq c d^{3/2} H^3 n^{-1} \log(dH/\delta), \quad (11a)$$

and is valid for a sample size $n \geq cd^2 H^3$. Alternatively stated, Corollary 1 guarantees that for FQI using ridge regression with d -dimensional function approximation, the number of samples $n(\epsilon)$ required to obtain ϵ -optimal policy is at most

$$n_{\text{fast}}(\epsilon) \asymp d^{3/2} H^3 / \epsilon + d^2 H^3, \quad (11b)$$

where we use \asymp to denote a scaling that ignores constants and logarithmic factors.

Let us compare this guarantee to related work by Zanette et al. [37], who analyzed the use of pessimistic actor-critic methods for linear function classes. When translated into the notation of our paper, their analysis established⁵ a sample complexity of the order $n_{\text{Zan}}(\epsilon) \asymp d^2 H^3 / \epsilon^2$. Consequently, we see that once the target error ϵ is relatively small— $\epsilon \in (0, 1)$ —then stable MDPs can exhibit a much smaller $(1/\epsilon)$ sample complexity.

It should be noted that past work (e.g., [21, 37]) has established $(1/\epsilon^2)$ -lower bounds on the sample complexity of estimating ϵ policies in the off-line setting. However, these lower bounds *do not* contradict our fast rate guarantee (11b), because the “hard instances” used in these lower bound proofs violate the stability condition (**Stb**(ξ)). In particular, even infinitesimally small perturbations in policy lead to occupation measures that are significantly different.

When is pessimism necessary? An interesting aspect of the guarantee from Corollary 1 is that it provides guarantees for off-policy RL (and with fast rates) using a method that does *not* incorporate any form of pessimism. This is a sharp contrast with many other methods for off-policy RL, such as pessimistic forms of Q -learning and actor-critic methods (e.g., [21, 37]).

To be clear, as noted following the bound (11a), the guarantee from Corollary 1 requires the sample size to be lower bounded as $n \geq cd^2 H^3$. In contrast, pessimistic schemes only require a sample size sufficiently large to ensure validity of the Bellman residual upper bounds that underlie Corollary 1—meaning that $n \gtrsim d$ up to logarithmic factors. Thus, the pessimism principle can be useful for problems with smaller sample sizes.

⁵See Appendix C.2 for the details of this calculation.

3.2 Consequences for on-line RL

In this section, we explore some consequences of Theorem 1 for on-line reinforcement learning. We begin by describing a two-stage procedure⁶ that allows us to convert the risk bounds for FQI from off-line RL into regret in on-line RL:

Phase 1 (Exploration) In the initial T_0 episodes, the focus is purely on exploration, resulting in an estimate of Q -function denoted as $\hat{f}^{(T_0)}$.

Phase 2 (Fine-tuning) For $k = 0, 1, \dots, K-1$ with $K := \lceil \log_2(T/T_0) \rceil$, repeat:

- In the t -th episode, for each $t = T_0 2^k + 1, \dots, T_0 2^{k+1}$, execute the greedy policy induced by function $\hat{f}^{(T_0 2^k)}$.
- Update the Q -function estimate $\hat{f}^{(T_0 2^{k+1})}$ using FQI based on observations collected from episodes $T_0 2^k + 1, T_0 2^k + 2, \dots, T_0 2^{k+1}$.

We assume the burn-in time T_0 is large enough so as to ensure the pilot Q -function estimate $\hat{f}^{(T_0)}$ obtained in Phase 1 falls within a certain “absorbing” region $\mathcal{N}(\rho)$ around Q^* . Under these conditions, we have the following bound on the regret.

Corollary 2. For FQI based on ridge regression with rewards in $[0, 1]$, with a sufficiently large burn-in time T_0 and with suitable choices of the regularization parameters $\{\lambda_h\}_{h=1}^{H-1}$, the two-phase scheme achieves regret bounded as

$$\text{Regret}(T) \leq c \{T_0 \cdot H + d\sqrt{d} H^4 \log T \cdot \log(dHK/\delta)\}$$

with probability at least $1 - \delta$.

See Appendix C.1 for the proof.

Sharper bound on regret: The leading term (as T grows) in the regret bound grows as $\log T$, which is much smaller than the typical \sqrt{T} -rate found in past work [18, 20]. The \sqrt{T} rate has been shown to be unimprovable in general, but the worst-case instances [18, 20] that lead to \sqrt{T} -regret violate the stability conditions used in our analysis.

When is optimism needed? The use of optimism—by adding bonuses to the current value function estimates so as to encourage exploration—underlies many schemes in on-line RL. An interesting take-away from Corollary 2 is that under the stability conditions highlighted by our theory, it is possible to achieve excellent regret bounds without the use of optimism. In our two-phase scheme, the only exploration occurs in Phase 1. All other data is simply collected using the greedy policy induced by the current Q -function estimate. A well-designed exploration scheme—one that might incorporate the optimism principle—is necessary only during the burn-in Phase 1.

4 Discussion

This paper introduces a novel approach for the analysis of value-based RL methods for continuous state-action spaces. Our analysis highlights two key stability properties of MDPs under which much sharper bounds on value sub-optimality can be guaranteed. Our analysis offers fresh perspectives on the commonly used pessimism and optimism principles, in off-line and on-line settings respectively.

Our study leaves open various questions for future work. First, our main result (Theorem 1) has consequences for linear quadratic control, to be described in an upcoming paper [8]. It provides insight into the role of covariate shift in linear quadratic control, as well as efficient exploration in the on-line setting. Second, our current statistical analysis focused on i.i.d. data with linear function approximation. It is interesting to consider the extensions to dependent data and non-parametric function approximation (e.g. kernels, boosting, and neural networks). Third, while this paper has provided upper bounds, it remains to address the complementary question of lower bounds for policy optimization over the classes of stable MDPs isolated here. Last, to better align our framework with real-world scenarios, we intend to go beyond the idealized completeness condition used in this paper, and treat the role of model mis-specification.

⁶To be clear, the purpose of this scheme is primarily conceptual, rather than practical in nature.

Acknowledgements

This work was partially supported by NSF grant CCF-1955450, ONR grant N00014-21-1-2842, and NSF DMS-2311072 to MJW.

References

- [1] A. Altamimi, C. Lagoa, J. G. Borges, M. E. McDill, C. Andriotis, and K. Papakonstantinou. Large-scale wildfire mitigation through deep reinforcement learning. *Frontiers in Forests and Global Change*, 5:734330, 2022.
- [2] H. Bastani, M. Bayati, and K. Khosravi. Mostly exploration-free algorithms for contextual bandits. *Management Science*, 67(3):1329–1349, 2021.
- [3] D. Bertsekas. *Lessons from AlphaZero for optimal, model predictive, and adaptive control*. Athena Scientific, 2022.
- [4] 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.
- [5] J. Degraeve, F. Felici, and J. B. et al. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602:414–419, 2022.
- [6] 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*, pages 2826–2836. PMLR, 2021.
- [7] Y. Duan, C. Jin, and Z. Li. Risk bounds and Rademacher complexity in batch reinforcement learning. In *International Conference on Machine Learning*, pages 2892–2902. PMLR, 2021.
- [8] Y. Duan and M. J. Wainwright. Covariate shift in linear quadratic control. *Manuscript*.
- [9] Y. Duan and M. J. Wainwright. Policy evaluation from a single path: Multi-step methods, mixing and mis-specification. *arXiv preprint arXiv:2211.03899*, 2022.
- [10] Y. Duan and M. Wang. Minimax-optimal off-policy evaluation with linear function approximation. In *International Conference on Machine Learning*, pages 2701–2709. PMLR, 2020.
- [11] Y. Duan, M. Wang, and M. J. Wainwright. Optimal policy evaluation using kernel-based temporal difference methods. *arXiv preprint arXiv:2109.12002*, 2021.
- [12] D. J. Foster, S. M. Kakade, J. Qian, and A. Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- [13] J. Gijbrecchts, R. N. Boute, J. A. Van Mieghem, and D. J. Zhang. Can deep reinforcement learning improve inventory management? Performance on lost sales, dual-sourcing, and multi-echelon problems. *Manufacturing & Service Operations Management*, 24(3):1349–1368, 2022.
- [14] A. Gonen and S. Shalev-Shwartz. Fast rates for empirical risk minimization of strict saddle problems. In *Conference on Learning Theory*, pages 1043–1063. PMLR, 2017.
- [15] E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69:169–192, 2007.
- [16] J. He, D. Zhou, and Q. Gu. Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pages 4171–4180. PMLR, 2021.
- [17] Y. Hu, N. Kallus, and M. Uehara. Fast rates for the regret of offline reinforcement learning. *Conference on Learning Theory*, 134:2462–2462, 2021.
- [18] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan. Is Q-learning provably efficient? *Advances in neural information processing systems*, 31, 2018.

- [19] C. Jin, Q. Liu, and S. Miryoosefi. Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418, 2021.
- [20] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- [21] Y. Jin, Z. Yang, and Z. Wang. Is pessimism provably efficient for offline RL? *International Conference on Machine Learning*, pages 5084–5096, 2021.
- [22] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. Al Sallab, S. Yogamani, and P. Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.
- [23] T. Koren and K. Levy. Fast rates for exp-concave empirical risk minimization. *Advances in Neural Information Processing Systems*, 28, 2015.
- [24] H. Mania, S. Tu, and B. Recht. Certainty equivalence is efficient for linear quadratic control. *Advances in Neural Information Processing Systems*, 32, 2019.
- [25] T. Nguyen-Tang, M. Yin, S. Gupta, S. Venkatesh, and R. Arora. On instance-dependent bounds for offline reinforcement learning with linear function approximation. *Association for the Advancement of Artificial Intelligence*, 2023.
- [26] M. L. Puterman and S. L. Brumelle. On the convergence of policy iteration in stationary dynamic programming. *Mathematics of Operations Research*, 4(1):60–69, 1979.
- [27] A. Rao and T. Jelvis. *Foundations of Reinforcement Learning with Applications to Finance*. CRC Press, Boca Raton, FL, 2022.
- [28] P. Rashidinejad, B. Zhu, C. Ma, J. Jiao, and S. Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, 34:11702–11716, 2021.
- [29] B. Recht. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems*, 2:253–279, 2019.
- [30] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354, 2017.
- [31] S. Spielberg, A. Tulsyan, N. P. Lawrence, P. D. Loewen, and R. B. Gopaluni. Toward self-driving processes: A deep reinforcement learning approach to control. *Amer. Inst. Chem. Eng. Journal*, 65:e16689, 2022.
- [32] L. Tai, G. Paolo, and M. Liu. Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 31–36. IEEE, 2017.
- [33] X. Wang, Q. Cui, and S. S. Du. On gap-dependent bounds for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 35:14865–14877, 2022.
- [34] T. Xie and N. Jiang. Q^* approximation schemes for batch reinforcement learning: A theoretical comparison. In *Conference on Uncertainty in Artificial Intelligence*, pages 550–559. PMLR, 2020.
- [35] M. Yin, Y. Duan, M. Wang, and Y.-X. Wang. Near-optimal offline reinforcement learning with linear representation: Leveraging variance information with pessimism. *arXiv preprint arXiv:2203.05804*, 2022.
- [36] C. Yu, J. Liu, S. Nemati, and G. Yin. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- [37] A. Zanette, M. J. Wainwright, and E. Brunskill. Provable benefits of actor-critic methods for offline reinforcement learning. *Advances in neural information processing systems*, 34:13626–13640, 2021.

A Proofs of Theorem 1

This section is devoted to the proof of Theorem 1, which consists of three main steps. These steps rely on two auxiliary lemmas whose proofs are fairly technical, so that they are deferred to in Appendices B.1 and B.2.

High-level outline: Let us outline the three steps of the proof. In Step 1, we use a one-step expansion of the difference in the occupation measures to reformulate the standard telescope inequality (6). Doing so results in a relation with structure similar to that of the left-hand side of inequality (Stb(ξ)). In Step 2, we develop a constraint on the function estimation error $d_h(\widehat{f}_h, Q_h^*)$ that ensures the occupation measure produced by policy $\widehat{\pi}$ remains stable and does not deviate too much from the occupation measure associated with the optimal policy π^* . In Step 3, we use Bellman stability (Stb(\mathcal{T})) to connect the Q -function error $\widehat{f}_h - Q_h^*$ with Bellman residuals. With this high-level view in place, we now work through the three steps.

A.1 Step 1: Reformulation of the telescope inequality.

Recall the standard telescope inequality (6). Our proof makes use of an alternative form, which involves the functions

$$\Delta_h(\pi; s, a) = \sum_{h'=h}^{H-1} \mathcal{P}_{h,h'}^\pi (\mathcal{T}_{h'}^* \widehat{f}_{h'+1} - \widehat{f}_{h'}) (s, a). \quad (12)$$

Lemma 1. *Given a Q -function estimate $\widehat{f} = (\widehat{f}_1, \dots, \widehat{f}_{H-1}, \widehat{f}_H = r_H)$ and the associated greedy policy $\widehat{\pi}$, we have the bound*

$$J(\pi) - J(\widehat{\pi}) \leq \sum_{h=1}^{H-1} \mathbb{E}_{\widehat{\pi}} [\Delta_h(\pi; s_h, \pi_h(s_h)) - \Delta_h(\pi; s_h, \widehat{\pi}_h(s_h))] \quad (13)$$

valid for any policy π .

See Appendix B.1 for the proof.

We apply the bound (13) with $\pi = \pi^*$. Following some algebra, we find that

$$J(\pi^*) - J(\widehat{\pi}) \leq \sum_{h=1}^{H-1} \sum_{h'=h}^{H-1} \widehat{\beta}(h, h') \cdot \varepsilon_{h'},$$

where $\varepsilon_{h'}$ is an upper bound on the Bellman residual $\|\mathcal{T}_{h'}^* \widehat{f}_{h'+1} - \widehat{f}_{h'}\|_{h'}$ as given in equation (4a). The term $\widehat{\beta}(h, h')$ is given by

$$\widehat{\beta}(h, h') := \sup_{f \in \partial \mathcal{F}: \|f\|_{h'} > 0} \left\{ \frac{1}{\|f\|_{h'}} \left| \mathbb{E}_{\widehat{\pi}} \left[(\mathcal{P}_{h,h'}^* f)(s_h, \pi_h^*(s_h)) - (\mathcal{P}_{h,h'}^* f)(s_h, \widehat{\pi}_h(s_h)) \right] \right| \right\}. \quad (14a)$$

We note that the left-hand side of inequality (Stb(ξ)) has a similar form to the term $\widehat{\beta}(h, h')$, differing only in that the expectation is taken over the occupation measure of running the optimal policy π^* , rather than the estimated policy $\widehat{\pi}$.

A.2 Step 2: Constraint to ensure stability

Our next step is to establish an upper bound on the coefficient $\widehat{\beta}(h, h')$ defined by the estimated policy $\widehat{\pi}$ in terms of the analogous quantity defined by the optimal policy π^* —namely, the coefficient

$$\beta(h, h') := \sup_{f \in \partial \mathcal{F}: \|f\|_{h'} > 0} \left\{ \frac{1}{\|f\|_{h'}} \left| \mathbb{E}_{\pi^*} \left[(\mathcal{P}_{h,h'}^* f)(s_h, \pi_h^*(s_h)) - (\mathcal{P}_{h,h'}^* f)(s_h, \widehat{\pi}_h(s_h)) \right] \right| \right\}. \quad (14b)$$

In order to do so, we demonstrate that a sufficiently small function estimation error $d_h(\widehat{f}_h, Q_h^*)$ ensures the inequality

$$\sum_{h=1}^{H-1} \sum_{h'=h}^{H-1} \widehat{\beta}(h, h') \cdot \varepsilon_{h'} \leq 2 \sum_{h=1}^{H-1} \sum_{h'=h}^{H-1} \beta(h, h') \cdot \varepsilon_{h'}. \quad (15)$$

Once we have established this bound, we can replace the term $\beta(h, h')$ with $\kappa_{h,h'}(\pi^*) \cdot \|\widehat{f}_h - Q_h^*\|_h / \|Q_h^*\|_h$, using the inequality **(Stb(ξ))**.

We summarize the result in the following auxiliary lemma:

Lemma 2. *Suppose that the function estimation errors satisfy $d_h(Q_h, Q_h^*) \leq \frac{1}{2b_{\mathcal{F}}}(H-h+1)^{-1}$ for $h = 2, 3, \dots, H-1$ and the sequence $\varepsilon = (\varepsilon_1, \dots, \varepsilon_{H-1}, \varepsilon_H = 0)$ satisfies the regularity condition (4b). Then we have*

$$J(\pi^*) - J(\widehat{\pi}) \leq 2 \sum_{h=1}^{H-1} \frac{\|\widehat{f}_h - Q_h^*\|_h}{\|Q_h^*\|_h} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\pi^*) \varepsilon_{h'} \right\}. \quad (16)$$

See Appendix B.2 for the proof.

A.3 Step 3: Connecting Q -function error and Bellman residuals

The remaining piece of the proof is to connect the function difference $\widehat{f}_h - Q_h^*$ with Bellman residuals $\mathcal{T}_h^* \widehat{f}_{h+1} - \widehat{f}_h$, using the stability condition **(Stb(\mathcal{T}))** on the Bellman operator \mathcal{T}^* . This is relatively straightforward: indeed, we claim that

$$\|\widehat{f}_h - Q_h^*\|_h \leq \sum_{h'=h}^{H-1} \kappa_{h,h'}(\mathcal{T}^*) \cdot \|\mathcal{T}_{h'}^* \widehat{f}_{h'+1} - \widehat{f}_{h'}\|_{h'}. \quad (17)$$

Recall that $Q_h^* = \mathcal{T}_h^* Q_{h+1}^*$ for $h = 1, 2, \dots, H-1$. Therefore, we have

$$\widehat{f}_h - Q_h^* = (\mathcal{T}_h^* \widehat{f}_{h+1} - \mathcal{T}_h^* Q_{h+1}^*) - (\mathcal{T}_h^* \widehat{f}_{h+1} - \widehat{f}_h).$$

By employing the triangle inequality and the Bellman stability given in equation **(Stb(\mathcal{T}))**, we derive that

$$\begin{aligned} \|\widehat{f}_h - Q_h^*\|_h &\leq \|\mathcal{T}_h^* \widehat{f}_{h+1} - \widehat{f}_h\|_h + \|\mathcal{T}_h^* \widehat{f}_{h+1} - \mathcal{T}_h^* Q_{h+1}^*\|_h \\ &\leq \|\mathcal{T}_h^* \widehat{f}_{h+1} - \widehat{f}_h\|_h + \kappa_h^* \|\widehat{f}_{h+1} - Q_{h+1}^*\|_{h+1}. \end{aligned}$$

Applying this inequality recursively yields the claim (17).

With this piece in place, we can complete the proof of Theorem 1. Indeed, we have

$$\begin{aligned} J(\pi^*) - J(\widehat{\pi}) &\stackrel{(a)}{\leq} 2 \sum_{h=1}^{H-1} \frac{\|\widehat{f}_h - Q_h^*\|_h}{\|Q_h^*\|_h} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\pi^*) \varepsilon_{h'} \right\} \\ &\stackrel{(b)}{\leq} 2 \sum_{h=1}^{H-1} \frac{1}{\|Q_h^*\|_h} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\mathcal{T}^*) \varepsilon_{h'} \right\} \left\{ \sum_{h'=h}^{H-1} \kappa_{h,h'}(\pi^*) \varepsilon_{h'} \right\}. \end{aligned}$$

Here step (a) is a restatement of the bound (16) from Lemma 2, whereas step (b) follows from inequality (17). Thus, we have established the claim given in Theorem 1.

B Proof of auxiliary lemmas for Theorem 1

We now turn to proofs of the two auxiliary results used to establish our main theorem, with Lemmas 1 and 2 treated in Appendices B.1 and B.2, respectively.

B.1 Proof of Lemma 1

For any integrable vector function $\mathbf{g} = (g_1, \dots, g_H) \in \mathbb{R}^{S \times A \times H}$, we define

$$D(\mathbf{g}) = \sum_{h=1}^H (\mathbb{E}_{\pi} - \mathbb{E}_{\hat{\pi}})[g_h(S_h, A_h)]. \quad (18a)$$

We claim that this functional satisfies the recursive relation

$$D(\mathbf{g}) = \sum_{h=1}^H \mathbb{E}_{\hat{\pi}}[g_h(S_h, \pi_h(S_h)) - g_h(S_h, \hat{\pi}_h(S_h))] + D(\mathcal{P}^{\pi} \mathbf{g}), \quad (18b)$$

where we have introduced the shorthand $\mathcal{P}^{\pi} \mathbf{g} := (\mathcal{P}_1^{\pi} g_2, \dots, \mathcal{P}_{H-1}^{\pi} g_H, 0) \in \mathbb{R}^{S \times A \times H}$.

Taking this claim as given for the moment, let us prove the bound (13) from Lemma 1. First, we set $\mathbf{g} := (\mathcal{P}^{\pi})^h \mathbf{g} = (\mathcal{P}_{1,1+h}^{\pi} g_{1+h}, \dots, \mathcal{P}_{H-h,H}^{\pi} g_H, 0, \dots, 0)$ in equation (18b) for $h = 0, 1, \dots, H-1$, which yields

$$D((\mathcal{P}^{\pi})^h \mathbf{g}) = \sum_{\substack{1 \leq h' \leq j \leq H, \\ j-h'=h}} \mathbb{E}_{\hat{\pi}}[\{\mathcal{P}_{h',j}^{\pi} g_j\}(S_{h'}, \pi_{h'}(S_{h'})) - \{\mathcal{P}_{h',j}^{\pi} g_j\}(S_{h'}, \hat{\pi}_{h'}(S_{h'}))] + D((\mathcal{P}^{\pi})^{h+1} \mathbf{g}).$$

Note that $(\mathcal{P}^{\pi})^H \mathbf{g} = 0$, which implies $D((\mathcal{P}^{\pi})^H \mathbf{g}) = 0$. We then sum the resulting bounds so as to obtain

$$D(\mathbf{g}) = \sum_{1 \leq h \leq h' \leq H} \mathbb{E}_{\hat{\pi}}[\{\mathcal{P}_{h,h'}^{\pi} g_{h'}\}(S_h, \pi_h(S_h)) - \{\mathcal{P}_{h,h'}^{\pi} g_{h'}\}(S_h, \hat{\pi}_h(S_h))]. \quad (19)$$

Setting $\mathbf{g} = \mathcal{T}^* \hat{\mathbf{f}} - \hat{\mathbf{f}}$, or equivalently $g_h = \mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h$, in equation (19), we find that

$$D(\mathcal{T}^* \hat{\mathbf{f}} - \hat{\mathbf{f}}) = \sum_{h=1}^{H-1} \mathbb{E}_{\hat{\pi}}[\Delta_h(\boldsymbol{\pi}; S_h, \pi_h(S_h)) - \Delta_h(\boldsymbol{\pi}; S_h, \hat{\pi}_h(S_h))],$$

where we have used the fact (12) that $\Delta_h(\boldsymbol{\pi}; \cdot) = \sum_{h'=h}^H \mathcal{P}_{h,h'}^{\pi} (\mathcal{T}_{h'}^* \hat{f}_{h'+1} - \hat{f}_{h'})$. Thus, we have established the bound (13) stated in Lemma 1.

It remains to establish the auxiliary claim (18b). Note that the functional D can be decomposed as $D(\mathbf{g}) = D_1 + D_2$, where

$$D_1 := \sum_{h=1}^H \mathbb{E}_{\hat{\pi}}[g_h(S_h, \pi_h(S_h)) - g_h(S_h, \hat{\pi}_h(S_h))] \quad \text{and}$$

$$D_2 := \sum_{h=1}^H (\mathbb{E}_{\pi} - \mathbb{E}_{\hat{\pi}})[g_h(S_h, \pi_h(S_h))].$$

Applying the tower property of conditional expectation, we find that

$$D_2 = \sum_{h=1}^{H-1} (\mathbb{E}_{\pi} - \mathbb{E}_{\hat{\pi}})[\mathbb{E}[g_{h+1}(S_{h+1}, \pi_{h+1}(S_{h+1})) \mid S_h, A_h]]$$

$$= \sum_{h=1}^{H-1} (\mathbb{E}_{\pi} - \mathbb{E}_{\hat{\pi}})[(\mathcal{P}_h^{\pi} g_{h+1})(S_h, A_h)] = D(\mathcal{P}^{\pi} \mathbf{g}).$$

Combining the expressions for D_1 and D_2 above yields the claim (18b).

B.2 Proof of Lemma 2

The key step in proving Lemma 2 is establishing that inequality (15) holds when the function estimation error $d_h(\widehat{f}_h, Q_h^*)$ is sufficiently small. In order to do so, we need to establish upper bounds on the term $\widehat{\beta}(h, h')$ by using $\beta(h, h')$. In particular, we will show that for any $1 \leq h \leq h' \leq H - 1$,

$$\widehat{\beta}(h, h') \leq \beta(h, h') + \sum_{j=1}^{h-1} \widehat{\beta}(j, h-1) \cdot b_{\mathcal{F}} \cdot d_h(\widehat{f}_h, Q_h^*). \quad (20)$$

The inequality (20) is derived based on the definitions of metric d_h and parameter $b_{\mathcal{F}} = 1$. After a close examination of the right-hand side of this inequality, it becomes evident that as long as the function estimation error $d_h(\widehat{f}_h, Q_h^*)$ is sufficiently small, the terms associated with $d_h(\widehat{f}_h, Q_h^*)$ are negligible and are dominated by $\beta(h, h')$. Consequently, inequality (15) within the arguments in Appendix A.2 is likely to hold true.

With claim (20) assumed to be valid at this point, we now establish a proper upper bound on the estimation error $d_h(\widehat{f}_h, Q_h^*)$ under which inequality (15) is satisfied. By taking linear combinations of inequality (20) using weights $\varepsilon = (\varepsilon_1, \dots, \varepsilon_{H-1}, \varepsilon_H = 0)$, we obtain

$$\begin{aligned} \sum_{h=1}^{H-1} \sum_{h'=h}^{H-1} \widehat{\beta}(h, h') \cdot \varepsilon_{h'} &\leq \sum_{h=1}^{H-1} \sum_{h'=h}^{H-1} \beta(h, h') \cdot \varepsilon_{h'} \\ &\quad + \sum_{h=2}^{H-1} \sum_{j=1}^{h-1} \widehat{\beta}(j, h-1) \cdot b_{\mathcal{F}} \cdot d_h(\widehat{f}_h, Q_h^*) \sum_{h'=h}^{H-1} \varepsilon_{h'}. \end{aligned} \quad (21)$$

When the sequence $\varepsilon = (\varepsilon_1, \dots, \varepsilon_{H-1}, \varepsilon_H = 0)$ is regular in the sense that inequality (4b) holds, the bound (21) reduces to

$$\begin{aligned} \sum_{1 \leq h \leq h' \leq H} \widehat{\beta}(h, h') \cdot \varepsilon_{h'} &\leq \sum_{1 \leq h \leq h' \leq H} \beta(h, h') \cdot \varepsilon_{h'} \\ &\quad + \sum_{1 \leq h \leq h' \leq H-2} \widehat{\beta}(h, h') \cdot \varepsilon_{h'} \cdot b_{\mathcal{F}} (H - h') \cdot d_{h'+1}(\widehat{f}_{h'+1}, Q_{h'+1}^*). \end{aligned}$$

Under the condition $d_h(\widehat{f}_h, Q_h^*) \leq \frac{1}{2b_{\mathcal{F}}}(H - h + 1)^{-1}$ for $2 \leq h \leq H - 1$, the inequality above implies bound (15), which further establishes the bound (16), as stated in Lemma 2.

It remains to prove the relation between $\widehat{\beta}(h, h')$ and $\beta(h, h')$, as shown in inequality (20).

Proof of bound (20): It is evident that inequality (20) holds for $h = 1$, therefore, we focus on its validation for indices $2 \leq h \leq H - 1$. Recall the definitions of functions $\widehat{\beta}(h, h')$ and $\beta(h, h')$, as given by equations (14a) and (14b). We apply the triangle inequality and derive that

$$\begin{aligned} &|\widehat{\beta}(h, h') - \beta(h, h')| \\ &\leq \sup_{f \in \partial \mathcal{F}: \|f\|_{h'} > 0} \left\{ \frac{1}{\|f\|_{h'}} \left| (\mathbb{E}_{\widehat{\pi}} - \mathbb{E}_{\pi^*}) \left[(\mathcal{P}_{h,h'}^* f)(S_h, \pi_h^*(S_h)) - (\mathcal{P}_{h,h'}^* f)(S_h, \widehat{\pi}_h(S_h)) \right] \right| \right\} \\ &= \sup_{f \in \partial \mathcal{F}: \|f\|_{h'} > 0} \left\{ \frac{1}{\|f\|_{h'}} \left| (\mathbb{E}_{\widehat{\pi}} - \mathbb{E}_{\pi^*}) \left[(\mathcal{P}_{h-1}^* - \mathcal{P}_{h-1}^{\widehat{\pi}}) \mathcal{P}_{h,h'}^* f \right](S_{h-1}, A_{h-1}) \right| \right\} \\ &=: \Delta\beta(h, h'). \end{aligned}$$

The term $\Delta\beta(h, h')$ involves differences from two sources: (i) the difference in transition kernels $\mathcal{P}_{h-1}^* - \mathcal{P}_{h-1}^{\widehat{\pi}}$ that captures the divergence between policies π_h^* and $\widehat{\pi}_h$; (ii) the discrepancy of occupation measures at the $(h - 1)$ -th step reflected by the difference in expectations $\mathbb{E}_{\pi^*} - \mathbb{E}_{\widehat{\pi}}$, which is determined by the policies $(\pi_1^*, \dots, \pi_{h-1}^*)$ and $(\widehat{\pi}_1, \dots, \widehat{\pi}_{h-1})$ until the $(h - 1)$ -th step. We treat them separately and write

$$\Delta\beta(h, h') \leq \nu_1(h - 1, h') \cdot \nu_2(h - 1), \quad (22)$$

where the functionals ν_2 and ν_1 are defined as

$$\begin{aligned}\nu_1(h-1, h') &:= \sup_{f \in \partial \mathcal{F}: \|f\|_{h'} > 0} \left\{ \frac{1}{\|f\|_{h'}} \left\| (\mathcal{P}_{h-1}^* - \mathcal{P}_{h-1}^{\widehat{\pi}}) \mathcal{P}_{h,h'}^* f \right\|_{h-1} \right\}, \\ \nu_2(h-1) &:= \sup_{f \in \partial \mathcal{F}: \|f\|_{h-1} > 0} \left\{ \frac{1}{\|f\|_{h-1}} \left| (\mathbb{E}_{\widehat{\pi}} - \mathbb{E}_{\pi^*}) [f(S_{h-1}, A_{h-1})] \right| \right\}.\end{aligned}$$

We first consider the term ν_1 . According to the definitions of metric d_h and parameter $b_{\mathcal{F}}$ we find that

$$\left\| (\mathcal{P}_{h-1}^* - \mathcal{P}_{h-1}^{\widehat{\pi}}) \mathcal{P}_{h,h'}^* f \right\|_{h-1} \leq d_h(\widehat{f}_h, Q_h^*) \cdot \left\| \mathcal{P}_{h,h'}^* f \right\|_h \stackrel{(*)}{\leq} d_h(\widehat{f}_h, Q_h^*) \cdot b_{\mathcal{F}} \|f\|_{h'},$$

which in turn implies

$$\nu_1(h-1, h') \leq b_{\mathcal{F}} \cdot d_h(\widehat{f}_h, Q_h^*). \quad (23a)$$

The proof of inequality (*) for $b_{\mathcal{F}} = 1$, as mentioned above, can be found in Appendix E.1.

As for term ν_2 , we claim that

$$\nu_2(h-1) \leq \sum_{j=1}^{h-1} \widehat{\beta}(j, h-1). \quad (23b)$$

Combining the bound $\widehat{\beta}(h, h') \leq \beta(h, h') + \Delta\beta(h, h')$ with inequalities (22), (23a) and (23b), we establish the bound (20), as claimed. It remains to prove the claim (23b).

Proof of inequality (23b): This proof is analogous to that of Lemma 1. We begin by introducing an analogue of the functional $D(g)$ from equation (18a); in particular, for any index $h \in [H-1]$ and function $g \in \partial \mathcal{F}$, define

$$D_h^*(g) := (\mathbb{E}_{\pi^*} - \mathbb{E}_{\widehat{\pi}}) [g(S_h, A_h)].$$

Using the notation of D_h^* , we can rewrite the left-hand side of inequality (23b) as $\nu_2(h-1) = \sup_{f \in \partial \mathcal{F}: \|f\|_{h-1} > 0} \{ |D_{h-1}^*(f)| / \|f\|_{h-1} \}$. Following the same arguments as in the proof of inequality (18b), we can show that

$$D_h^*(g) = \mathbb{E}_{\widehat{\pi}} [g(S_h, \pi_h^*(S_h)) - g(S_h, \widehat{\pi}_h(S_h))] + D_{h-1}^*(\mathcal{P}_{h-1}^* g) \quad \text{for } h = 1, 2, \dots, H, \quad (24)$$

where we set $D_0^* \equiv 0$.

We consider function $g := \mathcal{P}_{j,h-1}^* f$ for $1 \leq j < h \leq H-1$. It follows from equation (24) that

$$D_j^*(\mathcal{P}_{j,h-1}^* f) = \mathbb{E}_{\widehat{\pi}} [(\mathcal{P}_{j,h-1}^* f)(S_j, \pi_j^*(S_j)) - (\mathcal{P}_{j,h-1}^* f)(S_j, \widehat{\pi}_j(S_j))] + D_{j-1}^*(\mathcal{P}_{j-1,h-1}^* f),$$

where we have used the relation $\mathcal{P}_{j-1}^* \mathcal{P}_{j,h-1}^* = \mathcal{P}_{j-1,h-1}^*$. Recalling the definition of $\widehat{\beta}(j, h-1)$ in equation (14a), applying the triangle inequality yields

$$|D_j^*(\mathcal{P}_{j,h-1}^* f)| \leq \widehat{\beta}(j, h-1) \cdot \|f\|_{h-1} + |D_{j-1}^*(\mathcal{P}_{j-1,h-1}^* f)|.$$

Summing this equation over indices $j = 1, 2, 3, \dots, h-1$ yields

$$|D_{h-1}^*(f)| \leq \sum_{j=1}^{h-1} \widehat{\beta}(j, h-1) \cdot \|f\|_{h-1},$$

which establishes inequality (23b).

C Proof of corollaries

This section contains proofs of several corollaries.

C.1 Proof of Corollary 2

We now turn to proving Corollary 2 regarding ridge-based FQI in on-line settings.

In Phase 1 of pure exploration, the cumulative regret is always bounded from above by $T_0 \cdot H$. During Phase 2 of fine-tuning, we let $\hat{\pi}^k$ be the policy employed in the rounds $T_0 2^k + 1, T_0 2^k + 2, \dots, T_0 2^{k+1}$, which is determined by the estimate $\hat{f}^{(T_0 2^k)}$ calculated at the end of the $(T_0 2^k)$ -th round. To estimate the regret, we consider the decomposition

$$\sum_{t=T_0+1}^T \{J(\pi^*) - J(\hat{\pi}^{(t)})\} \leq \sum_{k=0}^{K-1} \sum_{t=T_0 2^k}^{T_0 2^{k+1}-1} \{J(\pi^*) - J(\hat{\pi}^{(t)})\} = \sum_{k=0}^{K-1} T_0 2^k \{J(\pi^*) - J(\hat{\pi}^k)\}.$$

We leverage our bound (11a) for off-line RL in Section 3.1.2 to control the value sub-optimality $J(\pi^*) - J(\hat{\pi}^k)$. Recall that the policy $\hat{\pi}^k$ is derived from i.i.d. trajectories collected from the rounds $T_0 2^{k-1} + 1, T_0 2^{k-1} + 2, \dots, T_0 2^k$. We divide those $T_0 2^{k-1}$ trajectories into $H - 1$ equal shares and use each share to conduct estimation in one iteration of the FQI procedure. This subsampling technique ensures the independence of samples used in different iterations. It is primarily adopted for the sake of convenience (to keep the explanations concise) and is not essential in general. It follows from inequality (11a) that the bound

$$J(\pi^*) - J(\hat{\pi}^k) \leq c \frac{d\sqrt{d} H^4}{T_0 2^k} \log(dHK/\delta)$$

holds uniformly for indices $k = 0, 1, \dots, K - 1$ with a probability exceeding $1 - \delta$.

Putting together the pieces, we arrive at

$$\text{Regret}(T) \leq T_0 \cdot H + c d\sqrt{d} H^4 K \log(dHK/\delta).$$

We then derive the regret bound in Corollary 2 by noticing that $K = \mathcal{O}(\log T)$.

C.2 Comparing to known off-line bounds

In this section, we derive the sample complexity $n_{\text{Zan}}(\epsilon) \asymp \frac{d^2 H^3}{\epsilon^2}$ in Section 3.1.2 based on the results of Zanette et al. [37]; it gives the conventional $1/\sqrt{n}$ slow rate to which we compare. Zanette et al. [37] proved upper bounds on a pessimistic actor-critic scheme based on d -dimensional linear function approximation. Using our notation, Theorem 1 in their paper [37] can be expressed as

$$J(\pi^*) - J(\hat{\pi}) \leq c \left\{ \frac{1}{H} \sum_{h=1}^{H-1} \sqrt{\bar{\phi}_h^\top (\hat{\Sigma}_{h,\mathcal{D}} + \lambda_h \mathbf{I})^{-1} \bar{\phi}_h} \right\} \sqrt{\frac{dH^4}{n}}, \quad (25)$$

where the vector $\bar{\phi}_h$ is given by $\bar{\phi}_h := \mathbb{E}_{\pi^*} [\phi(S_h, A_h)]$, the covariance matrix

$$\hat{\Sigma}_{h,\mathcal{D}} := \frac{1}{|\mathcal{D}_h|} \sum_{(s_{h,i}, a_{h,i}, s'_{h,i}, r_{h,i}) \in \mathcal{D}_h} \phi(s_{h,i}, a_{h,i}) \phi(s_{h,i}, a_{h,i})^\top.$$

We now consider the explicit dependence of this upper bound on dimension d , horizon H and sample size n . The divergence term $\bar{\phi}_h^\top (\hat{\Sigma}_{h,\mathcal{D}} + \lambda_h \mathbf{I})^{-1} \bar{\phi}_h$ measures the conditioning of the regularized covariance matrix $(\hat{\Sigma}_{h,\mathcal{D}} + \lambda_h \mathbf{I})$ along a specific direction of $\bar{\phi}_h$. When the feature mapping ϕ operates within a d -dimensional space, it is reasonable to assume that

$$\bar{\phi}_h^\top (\hat{\Sigma}_{h,\mathcal{D}} + \lambda_h \mathbf{I})^{-1} \bar{\phi}_h \leq c' d.$$

The bound (25) then reduces to $J(\pi^*) - J(\hat{\pi}) \leq cdH^2/\sqrt{n}$. Regarding the dependence on horizon H , we conjecture that by incorporating the law of total variance in a more refined manner, it may be possible to further reduce the dependence by a factor of \sqrt{H} . Under these conditions, the bound takes the form $J(\pi^*) - J(\hat{\pi}) \leq cd\sqrt{H^3/n}$.

D Details of the mountain car experiment

In this experiment, a car is situated in a valley between two hills. The car's objective is to overcome the gravitational pull and reach the top of the right hill by efficiently controlling its acceleration.

D.1 Structure of the Markov decision process

The Markov decision process underlying the mountain car problem has a state space $\mathcal{S} \subset \mathbb{R}^2$ and an action space $\mathcal{A} \subset \mathbb{R}$. The state $s = (p, v)$ consists of the current position p and velocity v , whereas the scalar action $a = f$ corresponds to the applied input force. The state variables (p, v) and action f are restricted as

$$\begin{aligned} p &\in [p_{\min}, p_{\max}] = [-1.2, 0.6], \\ v &\in [v_{\min}, v_{\max}] = [-0.07, 0.07] \quad \text{and} \\ f &\in [f_{\min}, f_{\max}] = [-1, 1]. \end{aligned}$$

The mountain is described by the function

$$m(p) = \frac{1}{3} \sin(3p) + \frac{0.025}{(p_{\max} - p)(p - p_{\min})},$$

over the interval $p \in [p_{\min}, p_{\max}]$.

Let m' be the derivative of the mountain shape function m , which represents the instantaneous slope, and let $(\sigma_v, \sigma_p) = (0.01, 0.0025)$ be a pair of standard deviations that dictate the amount of randomness in the updates. For an interval $[a, b]$, we define the truncation function

$$\Psi_{[a,b]}(u) := \begin{cases} u & \text{if } u \in [a, b], \\ b & \text{if } u > b, \\ a & \text{if } u < a. \end{cases}$$

With this notation, at each discrete time step $h = 0, 1, 2, \dots$, the position and velocity of the car evolve as

$$\begin{aligned} v_{h+1} &= \Psi_{[v_{\min}, v_{\max}]}(v_h + 0.0015 f_h - 0.0025 m'(p_h) + \sigma_v Z_h) \\ p_{h+1} &= \Psi_{[p_{\min}, p_{\max}]}(p_h + v_{h+1} + \sigma_p Z'_h) \end{aligned}$$

where (Z_h, Z'_h) are a pair of independent standard normal variables. Note that the system dynamics are non-linear due to both the presence of the derivative m' and the truncation function Ψ .

The objective of the car is to reach the peak of the mountain, designated by the position $p_{\text{goal}} = 0.45$. The reward at state-action pair (s, a) is given by

$$r(s, a) := -\frac{1}{10} f^2 + 100 [\max\{0, p - p_{\text{goal}}\}]^2.$$

For any policy π , we define the γ -discounted value function

$$J(\pi) := \mathbb{E}_{\pi} \left[\sum_{h=0}^{\infty} \gamma^h r(S_h, A_h) \right],$$

using $\gamma = 0.97$. The initial state $s_0 = (p_0, v_0)$ is generated with p_0 following a uniform distribution over the interval $[-0.6, -0.4]$, and we initialize with velocity $v_0 = 0$.

D.2 Fitted Q-iteration (FQI) with linear function approximation

Here we describe the use of fitted Q-iteration (FQI) with linear function approximation to estimate the optimal Q -function, along with the corresponding greedy policy $\hat{\pi}$.

Linear function approximation We approximate the the optimal Q -function $(s, a) \mapsto Q^*(s, a)$ using a d -dimensional linear function class with $d = 3000$ features. We begin by defining the *base*

feature maps $\phi_p : [p_{\min}, p_{\max}] \rightarrow \mathbb{R}^{50}$ for position, and $\phi_v : [v_{\min}, v_{\max}] \rightarrow \mathbb{R}^{15}$ for velocity, with components given by

$$\begin{cases} \phi_{p,2j+1}(p) := \cos(jp), & \text{for } j = 0, 1, \dots, 24, \text{ and} \\ \phi_{p,2j}(p) := \sin(jp), & \text{for } j = 1, 2, \dots, 25; \\ \phi_{v,2j+1}(v) := \cos(jv), & \text{for } j = 0, 1, \dots, 7, \text{ and} \\ \phi_{v,2j}(v) := \sin(jv), & \text{for } j = 1, 2, \dots, 7. \end{cases}$$

To represent the action $a \equiv f$, we define the *base action feature map*

$$\phi_f(f) := (1, f, f^2, f^3) \in \mathbb{R}^4.$$

The overall feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{3000}$ is constructed by taking the outer product of the three base feature maps ϕ_p , ϕ_v , and ϕ_f as follows:

$$\phi(s, a) := \text{vec}\{\phi_p(p) \otimes \phi_v(v) \otimes \phi_f(f)\} \in \mathbb{R}^{3000}. \quad (26)$$

Taking all possible triples of the three base features in the outer product leads to the overall dimension $d = 3000 = 50 \times 15 \times 4$. Given a weight vector $\mathbf{w} \in \mathbb{R}^{3000}$, we define the function $f_{\mathbf{w}}(s, a) := \langle \mathbf{w}, \phi(s, a) \rangle$, and we approximate the optimal Q -function using the function class $\mathcal{F} := \{f_{\mathbf{w}} \mid \mathbf{w} \in \mathbb{R}^{3000}\}$.

Fitted Q-iteration (FQI) We employed fitted Q-iteration with the linear feature $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{3000}$ to estimate an optimal policy $\hat{\pi}$. The FQI process begins by initializing the weight vector as $\mathbf{w}_0 := \mathbf{0} \in \mathbb{R}^{3000}$. In each iteration, we first use the dataset $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n \subset \mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S}$ to construct the pseudo-responses

$$y_i := r_i + \gamma \max_{a \in \mathcal{A}} \underbrace{\langle \mathbf{w}_t, \phi(s'_i, a) \rangle}_{f_{\mathbf{w}_t}(s'_i, a)} \quad \text{for } i = 1, \dots, n, \quad (27)$$

corresponding to a stochastic estimate of the Bellman update applied to our current Q -function estimate $f_{\mathbf{w}_t}$. The polynomial form of the force feature ϕ_f allows for a closed-form solution to the maximum operation required in equation (27). Given these pseudo-responses, we then update the weight vector $\mathbf{w}_t \rightarrow \mathbf{w}_{t+1}$ via the ridge regression

$$\mathbf{w}_{t+1} := \arg \min_{\mathbf{w} \in \mathbb{R}^{3000}} \left\{ \frac{1}{n} \sum_{i=1}^n \{y_i - \langle \mathbf{w}, \phi(s_i, a_i)\rangle\}^2 + \lambda_n \|\mathbf{w}\|_2^2 \right\}, \quad (28)$$

where $\lambda_n = \frac{0.01}{n}$ in all experiments reported here.

We terminate the procedure after at most 500 iterations, or when there have been 5 consecutive iterations with insignificant improvements in weights, where insignificant means that $\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2 / \sqrt{3000} < 0.005$. Letting $\hat{\mathbf{w}}$ represent the weight vector obtained from this procedure, the resulting policy $\hat{\pi}$ is given by selecting the greedy action based on the Q -function estimate $\hat{f}(s, a) := \langle \hat{\mathbf{w}}, \phi(s, a) \rangle$.

D.3 Experimental configurations

Our experiments were based on an off-line dataset consisting of n i.i.d. tuples

$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n \subset \mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S},$$

where the state-action pairs $\{(s_i, a_i) = (p_i, v_i, f_i)\}_{i=1}^n$ were generated from a uniform distribution over the cube $[p_{\min}, p_{\max}] \times [v_{\min}, v_{\max}] \times [f_{\min}, f_{\max}]$. We performed independent experiments with the sample size n varying over the range

$$\begin{aligned} n \in \{ \lfloor e^k \rfloor \mid k = 10.5, 10.75, 11, \dots, 13 \} \\ = \{36315, 46630, 59874, 76879, 98715, 126753, 162754, 208981, 268337, 344551, 442413\}. \end{aligned}$$

In each experiment, we generated a dataset \mathcal{D} , estimated an optimal policy $\hat{\pi}$ based on the data, and evaluated the return $J(\hat{\pi})$. For each sample size, we conducted 80 independent trials.

In order to evaluate the return $J(\hat{\pi})$, for each initial position $p_0 = -0.5 + 0.2j/1000$ with $j = -500, -499, -498, \dots, 499$, we simulated 30 independent 1000-step trajectories by executing the estimated policy $\hat{\pi}$. The average return over the 30×1000 trajectories is used as the estimate of $J(\hat{\pi})$.

In order to approximate the policy⁷ π^\dagger that represents “ground truth”, we conducted a single experiment with sample size $n = 6.4 \times 10^6$ to obtain π^\dagger . We simulated 1000 trajectories for each initial position p_0 and calculated the average return, which serves as the reference value $J(\pi^\dagger)$. The value sub-optimality is then computed as the difference $J(\pi^\dagger) - J(\hat{\pi})$.

E Verification of auxiliary claims

In this appendix, we collect the verification of various auxiliary claims made in the main text.

E.1 Properties of occupation measures

In this appendix, we prove a useful inequality

$$\|\mathcal{P}_{h,h'}^* f\|_h \leq \|f\|_{h'}, \quad (29)$$

which holds for the state-action occupation measures (3). This bound is used in the proof of bound (20) in Appendix B.2.

By definition, we have

$$\|\mathcal{P}_h^* f\|_h^2 = \mathbb{E}_{\pi^*} [(\mathcal{P}_h^* f)^2(S_h, A_h)] = \mathbb{E}_{\pi^*} \left[\mathbb{E}_h [f(S_{h+1}, \pi_{h+1}^*(S_{h+1})) \mid S_h, A_h]^2 \right].$$

According to the property of variance, we can deduce

$$\mathbb{E}_{\pi^*} \left[\mathbb{E}_h [f(S_{h+1}, \pi_{h+1}^*(S_{h+1})) \mid S_h, A_h]^2 \right] \leq \mathbb{E}_{\pi^*} [f^2(S_{h+1}, \pi_{h+1}^*(S_{h+1}))] = \|f\|_{h+1}^2.$$

As a consequence, we find that $\|\mathcal{P}_h^* f\|_h \leq \|f\|_{h+1}$. Applying this inequality recursively leads to the conclusion that for any indices $1 \leq h \leq h' \leq H$, we have

$$\|\mathcal{P}_{h,h'}^* f\|_h = \|\mathcal{P}_h^* \mathcal{P}_{h+1,h'}^* f\|_h \leq \|\mathcal{P}_{h+1,h'}^* f\|_{h+1} \leq \|\mathcal{P}_{h+2,h'}^* f\|_{h+2} \leq \dots \leq \|f\|_{h'}.$$

This establishes the bound (29).

E.2 Proof of the telescope inequality (6)

For completeness of this paper,⁸ let us prove the telescope relation (6) stated in Section 2.3. For any policy $\pi = (\pi_1, \dots, \pi_H)$ and sequence of functions $\mathbf{f} = (f_1, \dots, f_H)$ with $f_H = r_H$, we have the “telescope” relation

$$V_1^\pi(s) = f_1(s, \pi_1(s)) + \sum_{h=1}^{H-1} \mathbb{E}_\pi [(\mathcal{T}_h^\pi f_{h+1} - f_h)(S_h, A_h) \mid S_1 = s] \quad \text{for any state } s \in \mathcal{S}. \quad (30)$$

Here the value function V_1^π is given by $V_1^\pi(s) := Q_1^\pi(s, \pi_1(s))$. Taking $\mathbf{f} = \hat{\mathbf{f}}$ in equation (30) yields

$$V_1^\pi(s) = \hat{f}_1(s, \pi_1(s)) + \sum_{h=1}^{H-1} \mathbb{E}_\pi [(\mathcal{T}_h^\pi \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h) \mid S_1 = s]. \quad (31a)$$

Letting $\pi = \hat{\pi}$ in equation (31a) yields

$$V_1^{\hat{\pi}}(s) = \hat{f}_1(s, \hat{\pi}_1(s)) + \sum_{h=1}^{H-1} \mathbb{E}_{\hat{\pi}} [(\mathcal{T}_h^{\hat{\pi}} \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h) \mid S_1 = s]. \quad (31b)$$

⁷In general, it is not guaranteed that π^\dagger is equal to the optimal policy π^* , due to approximation error that might arise from using the linear function class defined here.

⁸We are not claiming novelty here; see Theorem 2 of the paper [34]; or Lemma 3.2 in the paper [7] for analogous results.

Since $\hat{\pi}$ is a greedy policy with respect to function \hat{f} , we have

$$\hat{f}_1(s, \hat{\pi}_1(s)) \geq \hat{f}_1(s, \pi_1(s)), \quad \text{and} \quad \mathcal{T}_h^{\hat{\pi}} \hat{f}_{h+1} = \mathcal{T}_h^* \hat{f}_{h+1} \geq \mathcal{T}_h^{\pi} \hat{f}_{h+1} \quad \text{for any policy } \pi.$$

Using this fact and subtracting equations (31a) and (31b), we obtain

$$V_1^{\pi}(s) - V_1^{\hat{\pi}}(s) \leq \sum_{h=1}^{H-1} (\mathbb{E}_{\pi} - \mathbb{E}_{\hat{\pi}}) [(\mathcal{T}_h^* \hat{f}_{h+1} - \hat{f}_h)(S_h, A_h) \mid S_1 = s].$$

Finally, taking the expectation over the initial distribution ξ_1 yields the claimed inequality (6).

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The paper thoroughly discusses each point made in the abstract, including fast-rate convergence of RL in continuous state-action spaces, key stability properties, and the pessimism and optimism principles.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Section 4 discusses the work's limitations and future directions, such as adapting the framework to linear quadratic control, extending beyond i.i.d. data and finite-dimensional linear function spaces, proving a lower bound to demonstrate sharpness, and examining model mis-specification.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The theoretical results in this paper are mathematically rigorous, with intuitions in the main body and detailed justifications and proofs in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Appendix D explains how to reproduce the synthetic Mountain Car experiment from Section 1.1, including MDP structure, linear function space construction, FQI implementation, and experimental setups. No data is required.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in

some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The submission provides a supplementary code document to reproduce the Mountain Car experiment.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Appendix D explains how to reproduce the synthetic Mountain Car experiment from Section 1.1, including MDP structure, linear function space construction, FQI implementation, and experimental setups.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Figure 1(b) in the paper includes error bars for the standard errors, and the estimated slope in Figure 1(b) uses a bootstrap confidence interval to confirm observation validity.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Section 1.1 details the computer setup: "The experiment ran for 3 days on two laptops, each equipped with an Apple M2 Pro CPU and 16 GB RAM."

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: Yes, the research in the paper fully conforms with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: The paper focuses on purely theoretical aspects of reinforcement learning and does not discuss direct societal impacts, given its theoretical nature.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: The paper does not use existing assets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.