

---

# Sample Complexity of Kernel-Based Q-Learning

---

Sing-Yuan Yeh<sup>1,2</sup> Fu-Chieh Chang<sup>1,3</sup> Chang-Wei Yueh<sup>4</sup> Pei-Yuan Wu<sup>3,4</sup> Alberto Bernacchia<sup>1</sup> Sattar Vakili<sup>1</sup>

<sup>1</sup>MediaTek Research

<sup>2</sup>Graduate Program of Data Science, National Taiwan University and Academia Sinica

<sup>3</sup>Graduate Institute of Communication Engineering, National Taiwan University

<sup>4</sup>Department of Electrical Engineering, National Taiwan University

## Abstract

Modern reinforcement learning (RL) often faces an enormous state-action space. Existing analytical results are typically for settings with a small number of state-actions, or simple models such as linearly modeled Q-functions. To derive statistically efficient RL policies handling large state-action spaces, with more general Q-functions, some recent works have considered nonlinear function approximation using kernel ridge regression. In this work, we derive sample complexities for kernel based Q-learning when a generative model exists. We propose a non-parametric Q-learning algorithm which finds an  $\epsilon$ -optimal policy in an arbitrarily large scale discounted MDP. The sample complexity of the proposed algorithm is order optimal with respect to  $\epsilon$  and the complexity of the kernel (in terms of its information gain). To the best of our knowledge, this is the first result showing a finite sample complexity under such a general model.

## 1 INTRODUCTION

In recent years, Reinforcement Learning (RL) has been successfully applied to several fields, including gaming (Silver et al., 2016; Lee et al., 2018; Vinyals et al., 2019), autonomous driving (Kahn et al., 2017), microchip design (Mirhoseini et al., 2021), robot control (Kalashnikov et al., 2018), and algorithm search (Fawzi et al., 2022). Real-world problems usually contain an enormous state-action space, possibly infinite. For example, the game of *Go* has  $10^{170}$  states (Silver et al., 2016), and the number of actions

in the space of algorithms for matrix multiplication is  $10^{30}$  (Fawzi et al., 2022). It is currently not fully understood how RL algorithms are able to learn successful policies to solve these problems. Modern function approximators, such as kernel-based learning and deep neural networks, seem to be required for this success.

An important theoretical question is as follows. Consider a Markov decision process (MDP) with an unknown transition probability distribution. Suppose that a generative model (Kakade, 2003) is available, which provides sample transitions from any state-action pair. How many samples are required to learn a sufficiently good policy? That is referred to as the sample complexity.

Previous works have derived theoretical bounds on the sample complexity, under certain simple settings such as tabular and linear MDPs. In the tabular setting, it was shown that the sample complexity of learning an  $\epsilon$ -optimal policy (that is the value function is at most  $\epsilon$  away from the optimal value function) is in  $\mathcal{O}(\frac{|\mathcal{S}||\mathcal{A}|}{\epsilon^2})$ , where  $|\mathcal{S}|$  and  $|\mathcal{A}|$  are the cardinality of the state and action sets, respectively (Kearns and Singh, 1998; Azar et al., 2013; Sidford et al., 2018a,b; Agarwal et al., 2020), implying that for a very large state-action space, a virtually infinite number of samples is required to obtain a good performance. Another line of work considers a linear MDP model, where the transition probability admits a linear representation in a  $d$ -dimensional state-action feature map. It was shown that the sample complexity is in  $\mathcal{O}(\frac{d}{\epsilon^2})$  in this case, that is independent of the size of state and action spaces (Yang and Wang, 2019). Unfortunately, the linear assumption is rather inflexible and not often the case in practice.

In order to address the limitations of small state-actions or simple models, arising from tabular and linear MDP assumptions, a few recent studies considered nonlinear function approximation over possibly infinite state-action domains using kernel ridge regression. In these works, the transition probability distribution (and sometimes the reward function) are flexibly represented using a kernel-

based model (Yang and Wang, 2020; Yang et al., 2020a,b). The kernel-based models provide powerful regressor and uncertainty estimates, which can be leveraged to guide the RL algorithm. Furthermore, kernel-based models have a great representation capacity and can model a wide range of problems, considering that all continuous functions on compact subsets of  $\mathbb{R}^d$  can be approximated using common kernels (Srinivas et al., 2010). The existing works, however, do not address the specific question of sample complexity considered in this work, and instead derive regret bounds in the setting of an episodic MDP. A more detailed comparison is provided in Section 1.2.

The kernel-based approaches may be understood as a linear model with an infinite-dimensional state-action feature map, that corresponds to, e.g., the Mercer eigenfeatures of the kernel (see Section 2.2). In this sense, the linear model is a special case of the kernel-based model with a linear kernel. Nonetheless, the results on the sample complexity of linear MDPs do not extend to the kernel-based models, as those sample complexities scale with the dimension of the feature map (that is possibly infinite in the kernel setting). In contrast, in the kernel setting, the sample complexity depends on certain kernel-specific properties determined by the complexity of the kernel, which will be discussed in more detail.

### 1.1 Contributions

Considering a discounted MDP and the question of sample complexity (similar to Azar et al., 2013; Sidford et al., 2018a,b; Yang and Wang, 2019), we extend and generalize the existing work as follows.

- We introduce Kernel-based Q-Learning, referred to as KQLearn, a sample collection algorithm, which returns an  $\epsilon$ -optimal policy with a finite sample complexity over a very general class of models. In comparison to tabular and linear MDP settings, KQLearn makes at least two innovative contributions. In the tabular setting, the samples are collected from all state-action pairs that leads to an  $|\mathcal{S}||\mathcal{A}|$  scaling of the sample complexity. In the linear setting, the samples are collected from a set of state-actions spanning the entire state-action space (leading to the scaling of the sample complexity with dimension  $d$ ). Then, an estimation of the parameters of the linear model are updated through value iteration. Neither of these approaches are feasible in our case with an infinite state-action space and a non-parametric kernel-based model. KQLearn instead takes advantage of uncertainties provided by the kernel model to create a finite state-action set which is used for collecting the samples. These samples are then passed through an approximate Bellman operator using kernel ridge regression to update the value function (that is a continuous

function over the entire state-action space).

- We derive a finite sample complexity for KQLearn under a wide range of kernel models. In particular, we consider two classes of kernels with exponentially ( $\sigma_m \sim \exp(-m^{\beta_e})$ ,  $\beta_e > 0$ ) and polynomially ( $\sigma_m \sim m^{-\beta_p}$ ,  $\beta_p > 1$ ) decaying Mercer eigenvalues  $\sigma_m$  (see Definition 2). We prove a sample complexity of  $\tilde{\mathcal{O}}(\frac{1}{\epsilon^2})$  and  $\tilde{\mathcal{O}}\left(\left(\frac{1}{\epsilon}\right)^{\frac{2\beta_p}{\beta_p-1}}\right)^1$  under these two settings, respectively. To the best of our knowledge, this is the first finite sample complexity, for all  $\beta_p > 1$ , and the first order optimal sample complexity in  $\epsilon$ , under the setting of polynomially decaying eigenvalues. Comparison with the related work is discussed in more detail in Section 1.2. As a special case, we recover the  $\tilde{\mathcal{O}}(\frac{d}{\epsilon^2})$  sample complexity of the linear setting reported in Yang and Wang (2019).

We acknowledge that our bounds on the sample complexity of KQLearn may not be order optimal in the dependence on the discount factor  $\gamma$ . In particular, our bounds grow with  $\frac{1}{(1-\gamma)^7}$  in the case of smooth kernels, similar to the PPQ-Learning algorithm proposed in Yang and Wang (2019) for the linear setting. Under the tabular and linear settings, however, this dependency was improved to  $\frac{1}{(1-\gamma)^3}$ , in Sidford et al. (2018a) and Yang and Wang (2019), respectively. It appears a challenging problem whether the same improvement is feasible here. As mentioned above, even establishing a finite sample complexity is a challenging problem, and the sample complexities in the existing work may diverge with difficult kernels (some polynomial kernels as discussed in Section 1.2).

### 1.2 Related Work

The specific problem of sample complexity in a discounted MDP using a generative model has been considered in tabular and linear settings. The results are summarized in Table 1. Other variants of the problem, consider MDPs in the absence of a generative model (e.g., see, Azar et al., 2017; Jin et al., 2018, 2020; Russo, 2019; Yang et al., 2020a,b; Kakade et al., 2020; Zhou et al., 2021; Domingues et al., 2021, as representative works, as well as references therein), often episodic, with  $T$  episodes of length  $H$ , and regret bounds depending on  $T$  and  $H$ . The regret bounds can then be translated into sample complexities (e.g., see, Jin et al., 2018; Yang et al., 2020b). These results are also reported in Table 1. Other approaches to nonlinear function approximation in RL include models with bounded *eluder* dimension (Wang et al., 2020; Ayoub et al., 2020) and smoothing kernels (Domingues et al., 2021). Among these works the two most relevant ones to ours are Yang and Wang (2019) and Yang et al. (2020a,b).

<sup>1</sup>The notations  $\mathcal{O}$  and  $\tilde{\mathcal{O}}$  are used to denote the mathematical order, and that up to hiding logarithmic factors, respectively.

Table 1: The existing sample complexities under various settings, discussed in Section 1.2.

ALGORITHM	MDP	SETTING	SAMPLE COMPLEXITY
(Q-learning with UCB, Jin et al., 2018)	Episodic	Tabular	$\tilde{O}\left(\frac{ \mathcal{S}  \mathcal{A} H^4}{\epsilon^2}\right)$
(LSVI-UCB, Jin et al., 2020)	Episodic	Linear	$\tilde{O}\left(\frac{d^\beta H^4}{\epsilon^2}\right)$
(KOVI, Yang et al., 2020b)	Episodic	Kernel-based, polynomial eigendecay	$\tilde{O}\left(H^4\left(\frac{1}{\epsilon}\right)^{\frac{2\beta_p}{\beta_p-2}}\right)$
		Kernel-based, exponential eigendecay	$\tilde{O}\left(\frac{H^4}{\epsilon^2}\right)$
(Sidford et al., 2018a, Variance-Reduced QVI)	Discounted	Tabular	$\tilde{O}\left(\frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^3\epsilon^2}\right)$
(PPQ-Learning, Yang and Wang, 2019)	Discounted	Linear	$\tilde{O}\left(\frac{d}{(1-\gamma)^7\epsilon^2}\right)$
(OPPQ-Learning, Yang and Wang, 2019)	Discounted	Linear	$\tilde{O}\left(\frac{d}{(1-\gamma)^3\epsilon^2}\right)$
<b>KQLearn</b>	Discounted	Kernel-based, polynomial eigendecay	$\tilde{O}\left(\frac{1}{\epsilon^{\frac{2\beta_p}{\beta_p-1}}(1-\gamma)^{\frac{7\beta_p-1}{\beta_p-1}}}\right)$
		Kernel-based, exponential eigendecay	$\tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^7}\right)$

Similar to Yang and Wang (2019), we also consider sample complexity in a discounted MDP using a generative model. We consider a non-parametric kernel-based model, while they considered a parametric linear model. Thus, neither their algorithm nor their results extend to our setting. The linear setting is a special case of the kernel setting with a linear kernel, in which, we recover the  $\tilde{O}\left(\frac{d}{\epsilon^2}\right)$  sample complexity, given in Yang and Wang (2019), for two algorithms: PPQ-Learning and OPPQ-Learning. The latter improved the sample complexity with respect to the discount factor.

Similar to Yang et al. (2020b), we also consider a kernel-based model. We consider sample complexity in a discounted MDP with a generative model, while they primarily considered regret bounds in an episodic MDP. They also reported sample complexities as a direct consequence of their regret bounds. Specifically, under the two settings of exponentially and polynomially decaying eigenvalues, their sample complexities translate to  $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$  and  $\tilde{O}\left(\left(\frac{1}{\epsilon}\right)^{\frac{2\beta_p}{\beta_p-2}}\right)$ , respectively. Under the polynomial setting, their sample complexity bound is larger than ours. In addition, their sample complexity is not always finite and may diverge when  $1 < \beta_p \leq 2$ , that includes many cases of interest. This suboptimality and possibly trivial result is a consequence of the superlinear (thus, trivial) regret bounds when  $1 < \beta_p \leq 2$ . See Vakili et al. (2021d), for a detailed discussion on the theoretical challenges related to this result.

For example, consider the Matérn family of kernels as one of the most commonly used (Snoek et al., 2012; Shahriari

et al., 2015) and theoretically interesting (Srinivas et al., 2010) family of kernels. For a Matérn kernel with smoothness parameter  $\nu$  on a  $d$  dimensional input domain,  $\beta_p = 1 + \frac{2\nu}{d}$  (Yang et al., 2020b). That implies the sample complexity in Yang et al. (2020b) diverges when  $d > 2\nu$  (that is often the case when using the Matérn kernel). We, however, emphasize that the discounted MDP with a generative model and the episodic MDPs are different settings, and cannot be compared directly. Nonetheless, we present the first always finite sample complexity under a very general setting covering all kernels with polynomially decaying eigenvalues.

Another related problem is the kernel-based bandit problem (Srinivas et al., 2010), which corresponds to a degenerate MDP with  $|\mathcal{S}| = 1$ . The kernel-based bandit problem is a well studied problem with order optimal regret bounds (Salgia et al., 2021; Li and Scarlett, 2022) and sample complexities (Vakili et al., 2021a). The lower bounds on sample complexities for the squared exponential (SE) and Matérn kernels are reported in (Scarlett et al., 2017), which have the same scaling with  $\epsilon$  (up to logarithmic factors) as in our results, showing the order optimality of our sample complexities with  $\epsilon$  (see Section 4). To obtain the optimal order of sample complexities, Vakili et al. (2021a) proposed the Maximum Variance Reduction (MVR) algorithm, which selects observation points with the highest uncertainty in a greedy manner. We build upon their approach by extending it to the sample complexities of Q-learning. This extension is significantly more involved due to the presence of Markovian dynamics and the Bellman operator.

There exist other studies in the field of kernel-based RL that are related to our work, including Ormoneit and Sen (2002); Domingues et al. (2021). However, there are notable differences between their work and ours in terms of both the settings and the algorithms. Specifically, these works assumed a Lipschitz reward function and a Lipschitz transition probability distribution, which allowed them to use kernel smoothing methods. In contrast, our work uses kernel ridge regression with the assumption of a bounded RKHS norm of the transition probability distribution.

**Paper structure:** In Section 2, the problem is formalized, after an overview of the background on MDPs and kernel ridge regression. In Section 3, KQLearn is presented. The results are discussed in Section 4. A high level analysis is provided in Section 5, while the details are deferred to the appendix.

## 2 PRELIMINARIES

In this section, we overview the background on MDPs and kernel ridge regression. We then formally state the problem of sample complexity for Q-learning under this setting.

### 2.1 Discounted Markov Decision Process

A discounted Markov Decision Process (MDP) can be described by the tuple  $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\gamma \in (0, 1)$  is the discount factor,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function and  $P(\cdot|s, a)$  is the transition probability distribution<sup>2</sup> on  $\mathcal{S}$  for the next state from state-action pair  $(s, a)$ . We use the notation  $\mathcal{Z} = \mathcal{S} \times \mathcal{A}$  to denote the state-action space. Our results generally hold true for (possibly very large and) finite  $\mathcal{Z}$  or certain infinite  $\mathcal{Z}$ . For correctness, we assume that  $\mathcal{Z}$  is a compact subset of  $\mathbb{R}^d$ .

The goal is to find a (possibly random) policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , that maximizes the long-term expected reward, i.e., the value function,

$$V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right],$$

where  $s_t \sim P(\cdot|s_{t-1}, \pi(s_{t-1}))$  forms the trajectory of the states. It can be shown that (e.g., see Puterman, 2014), under mild assumptions (e.g., continuity of  $P$ , compactness of  $\mathcal{Z}$ , and boundedness of  $r$ ) there exists an optimal policy  $\pi^*$  which attains the maximal possible value  $V^*$  at every state,

$$\forall s \in \mathcal{S} : V^*(s) := V^{\pi^*}(s) = \max_{\pi} V^\pi(s).$$

<sup>2</sup>We intentionally do not use the standard term *transition kernel* for  $P$ , to avoid confusion with the term *kernel* in kernel-based learning.

To simplify the notation, for a value function  $V : \mathcal{S} \rightarrow \mathbb{R}$ , let

$$[PV](s, a) := \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')].$$

The Q-function, also sometimes referred to as the state-action value function, of a policy  $\pi$ , and the optimal Q-function are defined as

$$Q^\pi(s, a) = r(s, a) + \gamma[PV^\pi](s, a), \text{ and} \\ Q^*(s, a) = Q^{\pi^*}(s, a),$$

respectively. The Bellman operator  $\mathcal{T} : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$  is defined as

$$\forall s \in \mathcal{S} : [\mathcal{T}V](s) = \max_{a \in \mathcal{A}} \{r(s, a) + \gamma[PV](s, a)\}.$$

**Sample complexity of  $\epsilon$ -optimal policies:** An  $\epsilon$ -optimal policy is defined as follows.

**Definition 1.** ( *$\epsilon$ -optimal policy*) A policy  $\pi$  is called  $\epsilon$ -optimal if it achieves near optimal values from any initial state as follows:

$$V^\pi(s) \geq V^*(s) - \epsilon, \quad \forall s \in \mathcal{S},$$

or equivalently  $\|V^\pi - V^*\|_\infty \leq \epsilon$ .

We aim to learn  $\epsilon$ -optimal policies using a small number of samples. In this work, following Kearns and Singh (1998); Azar et al. (2013); Sidford et al. (2018a,b); Yang and Wang (2019), we suppose that a generative model (Kakade, 2003) is given where the RL algorithm is able to query transition samples  $s' \sim P(\cdot|s, a)$  for any state-action pair  $(s, a) \in \mathcal{Z}$ . The *sample complexity* of an RL algorithm is defined as the number of such samples used by the algorithm to obtain an  $\epsilon$ -optimal policy.

### 2.2 RKHS and Kernel Ridge Regression

The existing work achieving finite sample complexity in the RL setting typically assumes a small state-action space or linearly modeled MDPs. These results can be generalized and extended using kernel-based learning. In particular, a natural approach is to use elements of a known reproducing kernel Hilbert space (RKHS) to model the transitions. In this section, we overview RKHSs and kernel ridge regression.

Let  $K : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  be a known positive definite kernel with respect to a finite Borel measure. Let  $\mathcal{H}_K$  be the RKHS induced by  $K$ , where  $\mathcal{H}_K$  contains a family of functions defined on  $\mathcal{Z}$ . Let  $\langle \cdot, \cdot \rangle_{\mathcal{H}_K} : \mathcal{H}_K \times \mathcal{H}_K \rightarrow \mathbb{R}$  and  $\|\cdot\|_{\mathcal{H}_K} : \mathcal{H}_K \rightarrow \mathbb{R}$  denote the inner product and the norm of  $\mathcal{H}_K$ , respectively. The reproducing property implies that for all  $f \in \mathcal{H}_K$ , and  $z \in \mathcal{Z}$ ,  $\langle f, K(\cdot, z) \rangle_{\mathcal{H}_K} = f(z)$ . Without loss of generality, we assume  $K(z, z) \leq 1$  for

all  $z$ . Mercer theorem implies, under certain mild conditions,  $K$  can be represented using an infinite dimensional feature map:

$$K(z, z') = \sum_{m=1}^{\infty} \sigma_m \psi_m(z) \psi_m(z'). \quad (1)$$

A formal statement and the details are provided in Appendix A.

**Kernel ridge regression:** Kernel-based models provide powerful regressor and uncertainty estimators (roughly speaking, surrogate posterior variances) which can be leveraged to guide the RL algorithm. In particular, consider an unknown function  $f \in \mathcal{H}_K$ . Consider a set  $\mathcal{U}_J = \{z_j\}_{j=1}^J \subset \mathcal{Z}$  of  $J$  inputs. Assume  $J$  noisy observations  $\{Y(z_j) = f(z_j) + \epsilon_j\}_{j=1}^J$  are provided, where  $\epsilon_j$  are i.i.d. zero mean sub-Gaussian noise terms. Kernel ridge regression provides the following regressor and uncertainty estimate, respectively (see, e.g., Schölkopf et al., 2002),

$$\begin{aligned} \hat{f}_{\mathcal{U}_J}(z) &= k_{\mathcal{U}_J}^\top(z) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}, \\ \Sigma_{\mathcal{U}_J}^2(z) &= K(z, z) - k_{\mathcal{U}_J}^\top(z) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} k_{\mathcal{U}_J}(z), \end{aligned} \quad (2)$$

where  $k_{\mathcal{U}_J}(z) = [K(z, z_1), \dots, K(z, z_J)]^\top$  is a  $J \times 1$  vector of the kernel values between  $z$  and observations,  $K_{\mathcal{U}_J} = [K(z_i, z_j)]_{i,j=1}^J$  is the  $J \times J$  kernel matrix,  $Y_{\mathcal{U}_J} = [Y(z_1), \dots, Y(z_J)]^\top$  is the  $J \times 1$  observation vector,  $I_J$  is the identity matrix of dimensions  $J$ , and  $\lambda > 0$  is a free regularization parameter.

**Confidence intervals:** The prediction and uncertainties provided by kernel ridge regression allow us to use standard confidence intervals in the algorithm and analysis. In particular, various results exist stating that with probability at least  $1 - \delta$ , the prediction function satisfies  $|f(z) - \hat{f}_{\mathcal{U}_J}(z)| \leq \beta(\delta) \Sigma_{\mathcal{U}_J}(z)$  (either for fixed  $z$ , or simultaneously for all  $z$ ) where the confidence interval width multiplier  $\beta(\delta)$  depends on the properties of the observation noise and the complexity of  $f$  in terms of its RKHS norm (Srinivas et al., 2010; Abbasi-Yadkori, 2013; Vakili et al., 2021a, 2022). If the domain  $\mathcal{Z}$  is finite, the uniform confidence bounds readily follow from a union bound over the confidence intervals for a fixed  $z$ . For continuous domains, a discretization argument is typically used considering the following continuity assumption.

**Assumption 1.** *For each  $n \in \mathbb{N}$ , there exists a discretization  $\mathbb{Z}$  of  $\mathcal{Z}$  such that, for any  $f \in \mathcal{H}_K$  with  $\|f\|_{\mathcal{H}_K} \leq C_K$ , we have  $f(z) - f([z]) \leq \frac{1}{n}$ , where  $[z] = \arg \min_{z' \in \mathbb{Z}} \|z' - z\|_{l^2}$  is the closest point in  $\mathbb{Z}$  to  $z$ , and  $|\mathbb{Z}| \leq c C_K^d n^d$ , where  $c$  is a constant independent of  $n$  and  $C_K$ .*

Assumption 1 is a technical and mild assumption that holds for typical kernels such as SE and Matérn with  $\nu >$

1 (Srinivas et al., 2010; Chowdhury and Gopalan, 2017; Vakili et al., 2021a).

In our analysis, we use the following confidence interval for the RKHS elements.

**Lemma 1** (Vakili et al. (2021a, 2022)). *Consider a fixed design of observation points where  $\mathcal{U}_J$  is independent of the observation noise. When the noise terms are sub-Gaussian with parameter  $R^3$  and  $\|f\|_{\mathcal{H}_K} \leq C_K$ , the following each hold uniformly in  $z \in \mathcal{Z}$ , with probability  $1 - \delta$ ,*

$$\begin{aligned} f(z) &\leq \hat{f}_{\mathcal{U}_J}(z) + \beta(\delta) \Sigma_{\mathcal{U}_J}(z), \\ f(z) &\geq \hat{f}_{\mathcal{U}_J}(z) - \beta(\delta) \Sigma_{\mathcal{U}_J}(z), \end{aligned} \quad (3)$$

where  $\beta(\delta) = \mathcal{O}\left(C_K + \frac{R}{\lambda} \sqrt{d \log\left(\frac{J C_K}{\delta}\right)}\right)$ .

**Maximal information gain:** It is useful for our analysis to define maximal information gain  $\Gamma_{K,\lambda}$ , that is a kernel specific complexity term. It allows us to bound the total uncertainty in the kernel model using results similar to elliptical potential lemma (Carpentier et al., 2020). In particular, let us define

$$\Gamma_{K,\lambda}(J) = \sup_{\mathcal{U} \subset \mathcal{Z}, |\mathcal{U}| \leq J} \frac{1}{2} \log \det \left( I_J + \frac{1}{\lambda^2} K_{\mathcal{U}} \right). \quad (4)$$

Then, we have the following.

**Lemma 2** (Srinivas et al. (2010)). *For any set  $\mathcal{U}_J \subset \mathcal{Z}$ , we have*

$$\sum_{j=1}^J \Sigma_{\mathcal{U}_{j-1}}^2(z_j) \leq \frac{2}{\log(1 + 1/\lambda^2)} \Gamma_{K,\lambda}(J). \quad (5)$$

### 2.3 Problem Formulation

Consider the discounted MDP described in Section 2.1. We are interested in designing an algorithm with a small sample complexity which obtains an  $\epsilon$ -optimal RL policy, under the assumption that the transition probability distribution lives in the RKHS of a known kernel. Without loss of generality, we assume its RKHS norm is bounded by 1.

**Assumption 2.** *Assume that the transition probability distribution satisfies,*

$$\forall s' \in \mathcal{S} : \|P(s' | \cdot, \cdot)\|_{\mathcal{H}_K} \leq 1.$$

This assumption is very flexible given the generality of the RKHSs. This is a standard assumption which is also used in Yang et al. (2020b). We do not make any explicit assumptions on the Q-function related to the policy. Recall the definition of  $PV$  from Section 2.1. In Lemma 3, we prove that for any  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ ,  $\|PV\|_{\mathcal{H}_K} \leq \frac{c}{1-\gamma}$ , as a consequence of Assumption 2, that is essential for our analysis.

<sup>3</sup>A random variable  $X$  is said to be sub-Gaussian with parameter  $R$  if its moment generating function is bounded by that of a zero mean Gaussian with variance  $R^2$ .

**Some generic notation:** For any real number  $x$ , and real numbers  $a, b$ , the notation  $\Pi_{[a,b]}[x]$  is used to denote the projection of  $x$  onto  $[a, b]$ . For any integer  $J$ ,  $I_J$  denotes the  $J \times J$  identity matrix, and  $\mathbf{0}_J$  denotes the  $J \times 1$  zero vector.

### 3 KERNEL BASED Q-LEARNING

In this section, we present a novel kernel based Q-learning algorithm, referred to as KQLearn. Recall  $Q(s, a) = r(s, a) + \gamma[PV](s, a)$ . The transition probability distribution  $P$  and the value function  $V$  are both unknown to the algorithm. The algorithm, thus, recursively approximates  $PV$ , in rounds, using kernel ridge regression of  $PV$  from the observations in the previous round. I.e., the algorithm performs updates based on an approximate Bellman operator using predictions for  $PV$  provided by the kernel model. The samples are collected based on uncertainties for  $PV$  in the kernel model. For this purpose, the algorithm first creates a maximum uncertainty set which is used to collect the samples.

**Maximum Uncertainty Set ( $\mathcal{U}_J$ ):** The algorithm starts with creating a maximum uncertainty set with size  $J$  referred to as  $\mathcal{U}_J \subset \mathcal{Z}$ . This set is created based on the uncertainties provided by the kernel model. In particular, each state-action is added to this set based on the following rule: choose the state-action with the highest uncertainty in the kernel model.

$$(s_j, a_j) = \arg \max_{(s,a) \in \mathcal{Z}} \Sigma_{\mathcal{U}_{j-1}}^2(s, a). \quad (6)$$

Then, recursively,  $\mathcal{U}_j = \mathcal{U}_{j-1} \cup \{(s_j, a_j)\}$ , starting from  $\mathcal{U}_0 = \emptyset$ . The set  $\mathcal{U}_J$  is then used to collect samples from the generative model.

The algorithm proceeds in rounds indexed by  $\ell = 1, \dots, L$ . Each round  $\ell$  receives noisy observations  $Y_{\mathcal{U}_j}^{(\ell-1)}$  of  $PV$  from the previous round,  $\ell - 1$ . These observations are then used within kernel ridge regression to form a regressor of  $PV$  over entire  $\mathcal{Z}$ , and obtain new observation  $Y^{(\ell)}$ . The observation vector is initialized to a zero vector  $Y^{(0)} = \mathbf{0}_J$ .

During each round  $\ell$ , for each state-action pair  $(s_j, a_j) \in \mathcal{U}_J$ , a transition state  $s'_j \sim P(\cdot | s_j, a_j)$  is acquired from the generative model. The observation  $Y^{(\ell)}(s_j, a_j)$  is then given as follows:

$$Y^{(\ell)}(s_j, a_j) = \Pi_{[0, \frac{1}{1-\gamma}]} \max_{a \in \mathcal{A}} \left\{ r(s'_j, a) + \gamma k_{\mathcal{U}_j}(s'_j, a)^\top (K_{\mathcal{U}_j} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_j}^{(\ell-1)} \right\}. \quad (7)$$

The second term on the right hand side is the regressor in kernel ridge regression on  $PV$ , using  $Y_{\mathcal{U}_j}^{(\ell-1)}$  as a vector

of observations. In the analysis, we show a high probability bound on the error of this regression. The vector  $Y_{\mathcal{U}_j}^{(\ell)} = [Y^{(\ell)}(s_1, a_1), \dots, Y^{(\ell)}(s_J, a_J)]^\top$  can be understood as updated noisy observations of  $PV$  which is passed to the next round,  $\ell + 1$ . By definition of the value function and the assumption of bounded rewards, it can be easily checked that  $0 \leq V^*(s) \leq \frac{1}{1-\gamma}$ , for all  $s \in \mathcal{S}$ . We thus project the value of  $PV$  on  $[0, \frac{1}{1-\gamma}]$  interval.

KQLearn collects  $N = JL$  samples in total.<sup>4</sup> A pseudo-code is provided in Algorithm 1.

After collecting all samples, the KQLearn algorithm returns an RL policy  $\pi$  which selects the actions based on the following proxy  $Q$ -function:

$$\widehat{Q}^{(L)}(s, a) = r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}^{(L)}. \quad (8)$$

Specifically, when state  $s$  is observed, the policy  $\pi$  selects the action  $\pi(s) = \arg \max_{a \in \mathcal{A}} \widehat{Q}^{(L)}(s, a)$ . The second term on the right hand side is the kernel ridge regression of  $PV$  using the observation in round  $L$  of the KQLearn algorithm.

---

#### Algorithm 1 Kernel-based Q-learning (KQLearn)

---

**Input** Discounted MDP with a generative model, kernel  $K$ , regularization parameter  $\lambda > 0$ , and  $N > 0$   
**Output**  $\widehat{Q}^{(L)} : \mathcal{Z} \rightarrow \mathbb{R}$

- 1: Initialize  $L, J \in \mathbb{N}, N = LJ$ .
  - 2: Initialize  $Y^{(0)} = \mathbf{0}_J$  and the set  $\mathcal{U}_0 = \emptyset$ .
  - 3: **for all**  $j = 1, \dots, J$  **do**
  - 4:     Update the function  $\Sigma_{\mathcal{U}_{j-1}}(\cdot)$  using Equation 2.
  - 5:     Pick  $(s_j, a_j) \leftarrow \arg \max_{(s,a) \in \mathcal{Z}} \Sigma_{\mathcal{U}_{j-1}}^2(s, a)$ .
  - 6:      $\mathcal{U}_j \leftarrow \mathcal{U}_{j-1} \cup \{(s_j, a_j)\}$ .
  - 7: **end for**
  - 8: **for all**  $\ell = 1, \dots, L$  **do** ▷ round
  - 9:     **for all**  $j = 1, \dots, J$  **do**
  - 10:         Obtain a sample transition state  $s'_j \sim P(\cdot | s_j, a_j)$
  - 11:         Update the  $Y^{(\ell)}$  as follows.  
 $Y^{(\ell)}(s_j, a_j) \leftarrow \Pi_{[0, \frac{1}{1-\gamma}]} \max_{a \in \mathcal{A}} \{ r(s'_j, a) + \gamma k_{\mathcal{U}_j}^\top(s'_j, a) (K_{\mathcal{U}_j} + \lambda I_J)^{-1} Y_{\mathcal{U}_j}^{(\ell-1)} \}$
  - 12:     **end for**
  - 13: **end for**
  - 14:  $\widehat{Q}^{(L)}(\cdot) = r(\cdot) + \gamma k_{\mathcal{U}_J}^\top(\cdot) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}^{(L)}$ .
- 

<sup>4</sup>For the simplicity of presentation, we assume  $N = JL$ . When  $J$  does not divide  $N$ , we can ignore the samples in the last round.

## 4 SAMPLE COMPLEXITY OF KQLEARN

In this section, we present our theoretical results. The following theorem establishes a bound on the error in the value function for the policy obtained in the KQLearn algorithm.

**Theorem 1.** *Consider the discounted MDP described in Section 2.1. Consider the KQLearn algorithm described in Section 3. Under Assumptions 1 and 2, with probability at least  $1 - \delta$ ,*

$$\|V^\pi - V^*\|_\infty \leq 2\beta(\delta) \left(\frac{\gamma}{1-\gamma}\right)^2 \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{J}} + 2\gamma^{L-1} \left(\frac{1}{1-\gamma}\right)^2,$$

where  $\beta(\delta) = \mathcal{O}\left(\frac{c}{1-\gamma} + \frac{1}{2\lambda(1-\gamma)} \sqrt{d \log\left(\frac{J}{(1-\gamma)\delta}\right)}\right)$  and  $c$  is a constant given in Lemma 3.

When  $J$  and  $L$  are sufficiently large, both terms in the upper bound become arbitrarily small.

### 4.1 Sample Complexity

We can obtain explicit sample complexities for the KQLearn algorithm, using kernel specific bounds on  $\Gamma_{K,\lambda}$ , which depend on the decay rate of the Mercer eigenvalues of  $K$ . In particular, we define the following characteristic eigendecay profiles (which are similar to those outlined in Chatterji et al. (2019); Vakili et al. (2021c); Yang et al. (2020b)).

**Definition 2.** *[Polynomial and Exponential Eigendecay] Consider the Mercer eigenvalues  $\{\sigma_m\}_{m=1}^\infty$  of  $K$  as given in Equation 1 in a decreasing order.*

- (i) *For some  $C_p > 0$ ,  $\beta_p > 1$ ,  $K$  is said to have a  $(C_p, \beta_p)$  polynomial eigendecay, if for all  $m \in \mathbb{N}$ , we have  $\sigma_m \leq C_p m^{-\beta_p}$ .*
- (ii) *For some  $C_{e,1}, C_{e,2}, \beta_e > 0$ ,  $K$  is said to have a  $(C_{e,1}, C_{e,2}, \beta_e)$  exponential eigendecay, if for all  $m \in \mathbb{N}$ , we have  $\sigma_m \leq C_{e,1} \exp(-C_{e,2} m^{\beta_e})$ .*

We are now ready to present explicit bounds on the sample complexity for the very general classes of kernels with polynomial and exponential decay of Mercer eigenvalues.

**Theorem 2.** *Consider the discounted MDP described in Section 2.1. Consider the KQLearn algorithm described in Section 3, with  $L = \Theta\left(\frac{\log(\epsilon(1-\gamma)^2)}{1-\gamma}\right)$  and  $J = \frac{N}{L}$ . Under Assumptions 1 and 2, KQLearn obtains an  $\epsilon$ -optimal policy with probability at least  $1 - \delta$ , with a sample complexity at most*

- *In the case of a kernel with  $(C_p, \beta_p)$  polynomial eigendecay,*

$$N = \tilde{\mathcal{O}}\left(\frac{\left(\log\left(\frac{1}{\delta}\right)\right)^{\frac{\beta_p}{\beta_p-1}}}{\epsilon^{\frac{2\beta_p}{\beta_p-1}} (1-\gamma)^{\frac{7\beta_p-1}{\beta_p-1}}}\right). \quad (9)$$

- *In the case of a kernel with  $(C_{e,1}, C_{e,2}, \beta_e)$  exponential eigendecay,*

$$N = \tilde{\mathcal{O}}\left(\frac{\log\left(\frac{1}{\delta}\right)}{\epsilon^2 (1-\gamma)^7}\right). \quad (10)$$

The detailed expressions for  $J$  and  $L$  including the implied constants and logarithmic factors in the  $\tilde{\mathcal{O}}$  notation are provided in Appendix B, Equations 29 and 32.

**Specific Kernels:** Our bounds on the sample complexity can be specialized for various kernels where the eigendecay or bounds on  $\Gamma_{K,\lambda}$  is known (such as the ones in Srinivas et al., 2010; Vakili et al., 2021c,a). Specifically, for the Matérn and SE kernels, we have, respectively,

$$N = \tilde{\mathcal{O}}\left(\frac{\left(\log\left(\frac{1}{\delta}\right)\right)^{1+\frac{d}{2\nu}}}{\epsilon^{2+\frac{d}{\nu}} (1-\gamma)^{7+\frac{3d}{\nu}}}\right), \text{ and } N = \tilde{\mathcal{O}}\left(\frac{\log\left(\frac{1}{\delta}\right)}{\epsilon^2 (1-\gamma)^7}\right).$$

### 4.2 Optimality of the Sample Complexities

The sample complexities given above are order optimal with respect to  $\epsilon$ . We compare them to the lower bounds on the sample complexity for kernel bandits (that is a special case of our setting when  $|S| = 1$ ). In particular, Scarlett et al. (2017) proved  $\Omega\left(\left(\frac{1}{\epsilon}\right)^{2+\frac{d}{\nu}}\right)$  and  $\Omega\left(\frac{1}{\epsilon^2}\right)$  sample complexities for the Matérn and SE kernels, respectively. Our results are the first finite sample complexities for the RL problem under a very general case which includes all kernels with polynomially decaying eigenvalues.

In terms of the discount factor, our sample complexities scale with  $\frac{1}{(1-\gamma)^7}$  in the case of smooth kernels (similar to the PPQ-Learning algorithm, in the linear setting, Yang and Wang, 2019). In the tabular and linear settings, however, this has been improved to  $\mathcal{O}\left(\frac{1}{(1-\gamma)^3}\right)$ . It remains an interesting problem for future investigation that whether the dependency of the sample complexity on the discount factor can be improved to  $\mathcal{O}\left(\frac{1}{(1-\gamma)^3}\right)$ , also in the kernel setting. As discussed in the introduction, in the kernel setting, neither observing all state-actions nor a parametric update of the model through value iteration is feasible. Thus, a different approach to algorithm design and analysis is required that is increasingly more challenging among these settings: *tabular*  $\rightarrow$  *linear*  $\rightarrow$  *kernel-based*.

## 5 ANALYSIS

Theorem 2 is a consequence of Theorem 1, and using the kernel specific bounds on  $\Gamma_{K,\lambda}$ . The proof of Theorem 1 builds on several components including tracking the approximation error in kernel ridge regression and convergence error of an approximate Bellman operator. In this section, we overview the main steps in the proof of Theorem 1, while deferring the details to the appendix.

**Approximate Bellman operator:** Recall the Bellman operator defined in Section 2.1. The transition probability distribution and the value function are complex non-linear functions on continuous domains, unknown to the algorithm. We thus define an approximate Bellman operator  $\widehat{\mathcal{T}}$ , which uses noisy observations of  $PV$  on a fixed set  $\mathcal{U}_J$ , and takes advantage of kernel ridge regression, to perform an approximate Bellman operator update. In particular, for all  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ , and a fixed set  $\mathcal{U}_J \subset \mathcal{Z}$ , let us define

$$[\widehat{\mathcal{T}}V](s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{PV}]_{\mathcal{U}_J} \right\}, \quad (11)$$

where  $[\widehat{PV}](s, a) = V(s')$  is a random variable,  $s' \sim P(\cdot | s, a)$  is a random transition state, and

$$[\widehat{PV}]_{\mathcal{U}_J} = \left[ [\widehat{PV}](s_1, a_1), \dots, [\widehat{PV}](s_J, a_J) \right]^\top.$$

In the KQLearn algorithm, define

$$\widehat{V}^{(\ell)}(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}^{(\ell)} \right\}, \quad (12)$$

which we refer to as proxy value function (similar to the proxy Q-function given in (8)). We then have the following recursive relation over  $\ell$ .

$$\widehat{V}^{(\ell)} = \widehat{\mathcal{T}}\Pi_{[0, \frac{1}{1-\gamma}]}[\widehat{V}^{(\ell-1)}]. \quad (13)$$

**Error in proxy value function:** In order to bound the error in the value function of the policy  $\pi$  obtained by KQLearn,  $\|V^\pi - V^*\|$ , we need to bound the error in the proxy Q-function given in (8), which is used to obtain  $\pi$ . The error in proxy Q-function can be bounded based on the error in the proxy value function at round  $L - 1$ . In particular, we have

$$\|\widehat{Q}^{(L)} - Q^*\|_\infty \leq \|\widehat{V}^{(L-1)} - V^*\|_\infty. \quad (14)$$

Therefore, we next bound the error in the proxy value function. We can write the error in proxy value function as the

sum of two terms: the error in approximate Bellman operator and the error in the value function using true Bellman operator. Specifically,

$$\begin{aligned} \|\widehat{V}^{(L-1)} - V^*\|_\infty &= \left\| \widehat{\mathcal{T}}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-2)} - V^* \right\|_\infty \\ &\leq \left\| \widehat{\mathcal{T}}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-2)} - \mathcal{T}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-2)} \right\|_\infty \\ &\quad + \left\| \mathcal{T}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-2)} - V^* \right\|_\infty. \end{aligned} \quad (15)$$

The second term can be recursively bounded which leads to the second term  $\frac{2\gamma^{L-1}}{(1-\gamma)^2}$  in the error bound in Theorem 1. The first term leads to an important step in the analysis which is based on the error in kernel ridge regression. Specifically, let  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$  be any value function. We have, for all  $s \in \mathcal{S}$ ,

$$\begin{aligned} &[\widehat{\mathcal{T}}V](s) - [\mathcal{T}V](s) \\ &= \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{PV}]_{\mathcal{U}_J} \right\} \\ &\quad - \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma [PV](s, a) \right\} \\ &\leq \gamma \max_{a \in \mathcal{A}} \left\{ k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{PV}]_{\mathcal{U}_J} \right. \\ &\quad \left. - [PV](s, a) \right\}. \end{aligned} \quad (16)$$

**Error in kernel ridge regression:** The term inside max in Equation 16 is the error in kernel ridge regression, where  $PV$  is the target function,  $[\widehat{PV}]_{\mathcal{U}_J}$  is a set of  $J$  noisy observations, and  $k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{PV}]_{\mathcal{U}_J}$  is the regressor. In order to apply Lemma 1, we need an upper bound on the RKHS norm of  $PV$ , as well as an upper bound on the sub-Gaussianity parameter of the observation noise in  $\widehat{PV}$ . These are established in the following lemmas.

**Lemma 3.** Consider an integrable value function  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ . Under Assumption 2, we have

$$\|PV\|_{\mathcal{H}_K} \leq \frac{c}{1-\gamma}, \quad (17)$$

where  $c = \int_{\mathcal{S}} ds$  is the volume of  $\mathcal{S}$ .

**Lemma 4.** Consider a transition probability distribution  $P$ , and an integrable value function  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ . We have, for all  $(s, a)$ ,  $\mathbb{E} \left[ [\widehat{PV}](s, a) \right] = PV(s, a)$ . In addition,  $[\widehat{PV}](s, a)$  is a sub-Gaussian random variable with parameter  $\frac{1}{2(1-\gamma)}$ .

Lemma 4 follows from the definition of  $\widehat{PV}$ , as well as Hoeffding lemma for bounded random variables. A proof of Lemma 3 is provided in Appendix C.



Applying Lemma 1, we obtain, with probability  $1 - \delta$ , for all  $(s, a) \in \mathcal{Z}$ ,

$$\left| k_{\mathcal{U}_J}^\top(s, a)(K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1}[\widehat{P}V]_{\mathcal{U}_J} - [PV](s, a) \right| \leq \beta(\delta)\Sigma_{\mathcal{U}_J}(s, a), \quad (18)$$

where  $\beta(\delta) = \mathcal{O}\left(\frac{c}{1-\gamma} + \frac{1}{2(1-\gamma)\lambda}\sqrt{d\log\left(\frac{Jc}{(1-\gamma)\delta}\right)}\right)$ .

Eventually, using Lemma 2 on the total uncertainty, and by the design of  $\mathcal{U}_J$ , we bound  $\Sigma_{\mathcal{U}_J}(s, a)$  on the right hand side.

We thus bounded the two terms in (21), that bounds the error in proxy value function. More details on the proof of theorems and the proof of lemmas are provided in Appendix B and Appendix C, respectively.

## 6 CONCLUSION

Modern RL often faces an enormous state-action space and complex models. We considered the question of sample complexity in a discounted MDP with a generative model under the kernel setting, furthering a line of research in the literature (e.g., see Kearns and Singh, 1998; Azar et al., 2017; Sidford et al., 2018a,b; Yang and Wang, 2019). We introduced a novel kernel-based Q learning algorithm referred to as KQLearn and proved a finite bound on its sample complexity for very general classes of kernels. That is to the best of our knowledge the first finite sample complexity result under the general kernel setting (including all kernels with polynomially decaying eigenvalues). In addition, compared to the lower bounds on the special case of the kernel bandit problem, our sample complexities are tight with respect to  $\epsilon$  in finding an  $\epsilon$ -optimal policy. Our sample complexities, however, scale possibly suboptimally with respect to the discount factor, which remains an interesting open problem for future investigation.

## References

- Abbasi-Yadkori, Y. (2013). Online learning for linearly parametrized control problems.
- Agarwal, A., Jiang, N., Kakade, S. M., and Sun, W. (2019). Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, pages 10–4.
- Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. (2020). Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474. PMLR.
- Azar, M. G., Munos, R., and Kappen, H. J. (2013). Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91(3):325–349.
- Azar, M. G., Osband, I., and Munos, R. (2017). Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR.
- Carpentier, A., Vernade, C., and Abbasi-Yadkori, Y. (2020). The elliptical potential lemma revisited. *arXiv preprint arXiv:2010.10182*.
- Chatterji, N., Pacchiano, A., and Bartlett, P. (2019). Online learning with kernel losses. In *Proceedings of Machine Learning Research*, volume 97, pages 971–980, Long Beach, California, USA. PMLR.
- Chowdhury, S. R. and Gopalan, A. (2017). On kernelized multi-armed bandits. In *International Conference on Machine Learning*, pages 844–853. PMLR.
- Christmann, A. and Steinwart, I. (2008). *Support Vector Machines*. Springer New York, NY.
- Domingues, O. D., Ménard, P., Pirotta, M., Kaufmann, E., and Valko, M. (2021). Kernel-based reinforcement learning: A finite-time analysis. In *International Conference on Machine Learning*, pages 2783–2792. PMLR.
- Fawzi, A., Balog, M., Huang, A., Hubert, T., Romera-Paredes, B., Barekatin, M., Novikov, A., R Ruiz, F. J., Schrittwieser, J., Swirszcz, G., et al. (2022). Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610(7930):47–53.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is q-learning provably efficient? *Advances in Neural Information Processing Systems*, 31.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR.
- Kahn, G., Villafior, A., Pong, V., Abbeel, P., and Levine, S. (2017). Uncertainty-aware reinforcement learning for collision avoidance. *arXiv preprint arXiv:1702.01182*.
- Kakade, S., Krishnamurthy, A., Lowrey, K., Ohnishi, M., and Sun, W. (2020). Information theoretic regret bounds for online nonlinear control. *Advances in Neural Information Processing Systems*, 33:15312–15325.
- Kakade, S. M. (2003). *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom).
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. (2018). Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673. PMLR.

- Kearns, M. and Singh, S. (1998). Finite-sample convergence rates for q-learning and indirect algorithms. In *Advances in Neural Information Processing Systems*, volume 11. MIT Press.
- Lee, K., Kim, S.-A., Choi, J., and Lee, S.-W. (2018). Deep reinforcement learning in continuous action spaces: a case study in the game of simulated curling. In *International Conference on Machine Learning*, pages 2937–2946. PMLR.
- Li, Z. and Scarlett, J. (2022). Gaussian process bandit optimization with few batches. In *International Conference on Artificial Intelligence and Statistics*, pages 92–107. PMLR.
- Mercer, J. (1909). Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 209:415–446.
- Mirhoseini, A., Goldie, A., Yazgan, M., Jiang, J. W., Songhori, E., Wang, S., Lee, Y.-J., Johnson, E., Pathak, O., Nazi, A., et al. (2021). A graph placement methodology for fast chip design. *Nature*, 594(7862):207–212.
- Ormoneit, D. and Sen, A. (2002). Kernel-based reinforcement learning. *Machine learning*, 49(2-3):161.
- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Russo, D. (2019). Worst-case regret bounds for exploration via randomized value functions. *Advances in Neural Information Processing Systems*, 32.
- Salgia, S., Vakili, S., and Zhao, Q. (2021). A domain-shrinking based bayesian optimization algorithm with order-optimal regret performance. *Advances in Neural Information Processing Systems*, 34:28836–28847.
- Scarlett, J., Bogunovic, I., and Cevher, V. (2017). Lower bounds on regret for noisy gaussian process bandit optimization. In *Conference on Learning Theory*, pages 1723–1742. PMLR.
- Schölkopf, B., Smola, A. J., Bach, F., et al. (2002). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press.
- Shahriari, B., Swersky, K., Wang, Z., Adams, R. P., and De Freitas, N. (2015). Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018a). Near-optimal time and sample complexities for solving markov decision processes with a generative model. *Advances in Neural Information Processing Systems*, 31.
- Sidford, A., Wang, M., Wu, X., and Ye, Y. (2018b). Variance reduced value iteration and faster algorithms for solving markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 770–787. SIAM.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 25.
- Srinivas, N., Krause, A., Kakade, S., and Seeger, M. (2010). Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML 2010 - Proceedings, 27th International Conference on Machine Learning*, pages 1015–1022.
- Vakili, S., Bouziani, N., Jalali, S., Bernacchia, A., and Shiu, D.-S. (2021a). Optimal order simple regret for gaussian process bandits. *Advances in Neural Information Processing Systems*, 34:21202–21215.
- Vakili, S., Bromberg, M., Garcia, J., Shiu, D.-s., and Bernacchia, A. (2021b). Uniform generalization bounds for overparameterized neural networks. *arXiv preprint arXiv:2109.06099*.
- Vakili, S., Khezeli, K., and Picheny, V. (2021c). On information gain and regret bounds in gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR.
- Vakili, S., Scarlett, J., and Javidi, T. (2021d). Open problem: Tight online confidence intervals for rkhs elements. In *Conference on Learning Theory*, pages 4647–4652. PMLR.
- Vakili, S., Scarlett, J., Shiu, D.-S., and Bernacchia, A. (2022). Improved convergence rates for sparse approximation methods in kernel-based learning. *arXiv preprint arXiv:2202.04005*.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. (2019). Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354.
- Wang, R., Salakhutdinov, R., and Yang, L. F. (2020). Provably efficient reinforcement learning with general value function approximation. *arXiv preprint arXiv:2005.10804*.
- Yang, L. and Wang, M. (2019). Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR.
- Yang, L. and Wang, M. (2020). Reinforcement learning in feature space: Matrix bandit, kernels, and regret

bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR.

Yang, Z., Jin, C., Wang, Z., Wang, M., and Jordan, M. (2020a). Provably efficient reinforcement learning with kernel and neural function approximations. *Advances in Neural Information Processing Systems*, 33:13903–13916.

Yang, Z., Jin, C., Wang, Z., Wang, M., and Jordan, M. I. (2020b). On function approximation in reinforcement learning: Optimism in the face of large state spaces. *arXiv preprint arXiv:2011.04622*.

Zhou, D., He, J., and Gu, Q. (2021). Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802. PMLR.

In the appendix, we provide some details and proofs omitted from the main paper due to space limit. In Appendix A, we provide a formal statement of Mercer theorem and a constructive definition of the RKHS. The proof of Theorems 1 and 2, and auxiliary lemmas are provided in Appendix B and Appendix C, respectively.

## A MERCER THEOREM

Mercer theorem (Mercer, 1909) provides a representation of the kernel in terms of an infinite dimensional feature map (see, e.g. Christmann and Steinwart (2008), Theorem 4.49). Let  $\mathcal{Z}$  be a compact metric space and  $\mu$  be a finite Borel measure on  $\mathcal{Z}$  (we consider Lebesgue measure in a Euclidean space). Let  $L^2_\mu(\mathcal{Z})$  be the set of square-integrable functions on  $\mathcal{Z}$  with respect to  $\mu$ . We further say a kernel is square-integrable if

$$\int_{\mathcal{Z}} \int_{\mathcal{Z}} K(z, z')^2 d\mu(z) d\mu(z') < \infty.$$

**Theorem 3. (Mercer Theorem)** *Let  $\mathcal{Z}$  be a compact metric space and  $\mu$  be a finite Borel measure on  $\mathcal{Z}$ . Let  $K$  be a continuous and square-integrable kernel, inducing an integral operator  $T_K : L^2_\mu(\mathcal{Z}) \rightarrow L^2_\mu(\mathcal{Z})$  defined by*

$$(T_K f)(\cdot) = \int_{\mathcal{Z}} K(\cdot, z') f(z') d\mu(z'),$$

where  $f \in L^2_\mu(\mathcal{Z})$ . Then, there exists a sequence of eigenvalue-eigenfunction pairs  $\{(\sigma_m, \psi_m)\}_{m=1}^\infty$  such that  $\sigma_m > 0$ , and  $T_K \psi_m = \sigma_m \psi_m$ , for  $m \geq 1$ . Moreover, the kernel function can be represented as

$$K(z, z') = \sum_{m=1}^{\infty} \sigma_m \psi_m(z) \psi_m(z'),$$

where the convergence of the series holds uniformly on  $\mathcal{Z} \times \mathcal{Z}$ .

According to Mercer representation theorem (see, e.g., Christmann and Steinwart (2008), Theorem 4.51), the RKHS induced by  $K$  can consequently be represented in terms of  $\{(\sigma_m, \psi_m)\}_{m=1}^\infty$ .

**Theorem 4. (Mercer Representation Theorem)** *Let  $\{(\sigma_m, \psi_m)\}_{i=1}^\infty$  be the Mercer eigenvalue eigenfunction pairs. Then, the RKHS of  $K$  is given by*

$$\mathcal{H}_K = \left\{ f(\cdot) = \sum_{i=1}^{\infty} w_i \sigma_i^{\frac{1}{2}} \psi_i(\cdot) : w_i \in \mathbb{R}, \|f\|_{\mathcal{H}_K}^2 := \sum_{i=1}^{\infty} w_i^2 < \infty \right\}.$$

Mercer representation theorem indicates that the scaled eigenfunctions  $\{\sqrt{\sigma_i} \psi_i\}_{i=1}^\infty$  form an orthonormal basis for  $\mathcal{H}_K$ .

## B PROOF OF THEOREMS

In this section, we provide the proof of main theorems.

### B.1 Proof of Theorem 1.

The proof of Theorem 1 builds on an approximate Bellman operator, that uses noisy observations of  $PV$  within the rounds of KQLearn, and kernel ridge regression. We prove bounds on the error of this approximate Bellman operator. That is then used to bound the error in the value function of the policy obtained by KQLearn.

**Approximate Bellman operator:** Recall the Bellman operator defined in Section 2.1. The transition probability distribution and the value function are complex non-linear functions on continuous domains, unknown to the algorithm. We thus define an approximate Bellman operator  $\widehat{T}$ , which uses noisy observations of  $PV$  on a fixed set  $\mathcal{U}_J$ , and takes advantage of kernel ridge regression, to perform an approximate Bellman operator update. In particular, for all  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ , and a fixed set  $\mathcal{U}_J \subset \mathcal{Z}$ , let us define

$$[\widehat{T}V](s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}V]_{\mathcal{U}_J} \right\},$$

where  $[\widehat{P}V](s, a) = V(s')$  is a random variable,  $s' \sim P(\cdot|s, a)$  is a random transition state, and

$$[\widehat{P}V]_{\mathcal{U}_J} = \left[ [\widehat{P}V](s_1, a_1), \dots, [\widehat{P}V](s_J, a_J) \right]^\top.$$

In the KQLearn algorithm, define

$$\widehat{V}^{(\ell)}(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}^{(\ell)} \right\},$$

which we refer to as proxy value function (similar to the proxy Q-function given in (8)). We then have the following recursive relation over  $\ell$ .

$$\widehat{V}^{(\ell)} = \widehat{\mathcal{T}}_{\Pi_{[0, \frac{1}{1-\gamma}]}}[\widehat{V}^{(\ell-1)}]. \quad (19)$$

**Error in proxy value function:** We next bound the error in the proxy value function. We can write the error in proxy value function as the sum of two terms: the error in approximate Bellman operator and the error in the value function using true Bellman operator. Specifically, for  $l > 1$ ,

$$\begin{aligned} \left\| \widehat{V}^{(l-1)} - V^* \right\|_\infty &= \left\| \widehat{\mathcal{T}}_{\Pi_{[0, \frac{1}{1-\gamma}]}} \widehat{V}^{(l-2)} - V^* \right\|_\infty \\ &\leq \underbrace{\left\| \widehat{\mathcal{T}}_{\Pi_{[0, \frac{1}{1-\gamma}]}} \widehat{V}^{(l-2)} - \mathcal{T}_{\Pi_{[0, \frac{1}{1-\gamma}]}} \widehat{V}^{(l-2)} \right\|_\infty}_{\text{Term I}} \\ &\quad + \underbrace{\left\| \mathcal{T}_{\Pi_{[0, \frac{1}{1-\gamma}]}} \widehat{V}^{(l-2)} - V^* \right\|_\infty}_{\text{Term II}}. \end{aligned} \quad (20)$$

We now bound the two terms on the right hand side of (20).

**Term I:** The first term leads us to an important step in the analysis which is based on the error in kernel ridge regression. Specifically, let  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$  be any value function. We have, for all  $s \in \mathcal{S}$ ,

$$\begin{aligned} [\widehat{\mathcal{T}V}](s) - [\mathcal{T}V](s) &= \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}V]_{\mathcal{U}_J} \right\} - \max_{a \in \mathcal{A}} \{ r(s, a) + \gamma [PV](s, a) \} \\ &\leq \gamma \max_{a \in \mathcal{A}} \left\{ k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}V]_{\mathcal{U}_J} - [PV](s, a) \right\}. \end{aligned} \quad (21)$$

**Error in kernel ridge regression:** The term inside max in Equation 21 is the error in kernel ridge regression, where  $PV$  is the target function,  $[\widehat{P}V]_{\mathcal{U}_J}$  is a set of  $J$  noisy observations, and  $k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}V]_{\mathcal{U}_J}$  is the regressor. In order to apply Lemma 1, we need an upper bound on the RKHS norm of  $PV$ , as well as an upper bound on the sub-Gaussianity parameter of the observation noise in  $\widehat{P}V$ . These are established in Lemmas 3 and 4, respectively. Specifically, we have

$$\|PV\|_{\mathcal{H}_K} \leq \frac{c}{1-\gamma}. \quad (22)$$

And,  $[\widehat{P}V](s, a)$  is a sub-Gaussian random variable with parameter  $\frac{1}{2(1-\gamma)}$ .

Applying Lemma 1, we obtain, with probability  $1 - \delta$ , for all  $(s, a) \in \mathcal{Z}$ ,

$$\left| k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}V]_{\mathcal{U}_J} - [PV](s, a) \right| \leq \beta(\delta) \Sigma_{\mathcal{U}_J}(s, a), \quad (23)$$

where  $\beta(\delta) = \mathcal{O} \left( \frac{c}{1-\gamma} + \frac{1}{2(1-\gamma)\lambda} \sqrt{d \log \left( \frac{Jc}{(1-\gamma)\delta} \right)} \right)$ .

**Bounding  $\Sigma_{\mathcal{U}_J}(s, a)$ :** Conditioning on a smaller subset of observation reduces the variance  $\Sigma_{\mathcal{U}_j}(s, a) \geq \Sigma_{\mathcal{U}_J}(s, a)$ , for all  $j \leq J$  (due to positive definiteness of the kernel matrix). By the selection rule of the observation points:

$$(s_j, a_j) = \arg \max_{(s, a) \in \mathcal{Z}} \Sigma_{\mathcal{U}_{j-1}}^2(s, a), \quad (24)$$

we have  $\Sigma_{\mathcal{U}_{j-1}}(s, a) \leq \Sigma_{\mathcal{U}_{j-1}}(s_j, a_j)$ . Thus, for all  $(s, a) \in \mathcal{Z}$ ,

$$\begin{aligned} \Sigma_{\mathcal{U}_J}^2(s, a) &\leq \frac{1}{J} \sum_{j=1}^J \Sigma_{\mathcal{U}_{j-1}}^2(s, a) \\ &\leq \frac{1}{J} \sum_{j=1}^J \Sigma_{\mathcal{U}_{j-1}}^2(s_j, a_j) \\ &\leq \frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}, \end{aligned}$$

where the last line follows from Lemma 2.

Replacing the bound on  $\Sigma_{\mathcal{U}_J}(s, a)$ , we obtain, for all  $V : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ ,  $s \in \mathcal{S}$ ,

$$[\widehat{\mathcal{T}}V](s) - [\mathcal{T}V](s) \leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}}. \quad (25)$$

Thus,

$$\left\| \widehat{\mathcal{T}}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(l-2)} - \mathcal{T}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(l-2)} \right\|_{\infty} \leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}}. \quad (26)$$

**Term II:** We now bound the second term on the right hand side of (20), by the contraction of Bellman operator.

$$\begin{aligned} \left\| \mathcal{T}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(l-2)} - V^* \right\|_{\infty} &\leq \gamma \left\| \Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(l-2)} - V^* \right\|_{\infty} \\ &\leq \gamma \left\| \widehat{V}^{(l-2)} - V^* \right\|_{\infty}. \end{aligned} \quad (27)$$

The second inequality follows from the observation that  $V^*(s) \in [0, \frac{1}{1-\gamma}]$  for all  $s \in \mathcal{S}$ .

Combing the bounds on Term I and Term II, we obtain

$$\left\| \widehat{V}^{(l-1)} - V^* \right\|_{\infty} \leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \gamma \left\| \widehat{V}^{(l-2)} - V^* \right\|_{\infty}.$$

Recursively bounding the error in proxy value function at round  $l$  using the error at round  $l - 1$  for  $l = 2, \dots, L - 1$ , we have,

$$\begin{aligned} \left\| \widehat{V}^{(L-1)} - V^* \right\|_{\infty} &\leq \beta(\delta) \sqrt{\frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}} \left( \sum_{i=1}^{L-1} \gamma^i \right) + \gamma^{L-1} \left\| \widehat{V}^{(0)} - V^* \right\|_{\infty} \\ &\leq \beta(\delta) \sqrt{\frac{2\Gamma_{K, \lambda}(J)}{\log(1 + 1/\lambda^2)J}} \left( \sum_{i=1}^{L-1} \gamma^i \right) + \frac{\gamma^{L-1}}{1-\gamma}, \end{aligned}$$

where the second inequality comes from  $\left\| \widehat{V}^{(0)} - V^* \right\|_{\infty} \leq \frac{1}{1-\gamma}$ .

Recall the definition of the proxy Q-function  $\widehat{Q}^{(L)}$  given in (8). We bound the error in  $\widehat{Q}^{(L)}$  as follows. For all  $(s, a) \in \mathcal{Z}$ ,

$$\begin{aligned}
 \left\| \widehat{Q}^{(L)}(s, a) - Q^*(s, a) \right\|_\infty &= \left\| \left( r(s, a) + \gamma k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} Y_{\mathcal{U}_J}^{(L)} \right) - \left( r(s, a) + \gamma [PV^*](s, a) \right) \right\|_\infty \\
 &\leq \left\| \gamma \left( k_{\mathcal{U}_J}^\top(s, a) (K_{\mathcal{U}_J} + \lambda^2 I_J)^{-1} [\widehat{P}\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)}]_{\mathcal{U}_J} - [P\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)}](s, a) \right) \right\|_\infty \\
 &\quad + \left\| \gamma \left( [P\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)}](s, a) - [PV^*](s, a) \right) \right\|_\infty \\
 &\leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \gamma \left\| \mathbb{E}_{s' \sim P(\cdot|s,a)} [\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)}] - \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^*] \right\|_\infty \\
 &\leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \gamma \left\| \Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)} - V^* \right\|_\infty \\
 &\leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \gamma \left\| \widehat{V}^{(L-1)} - V^* \right\|_\infty \\
 &\leq \gamma\beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \gamma \left( \beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} \left( \sum_{i=1}^{L-1} \gamma^i \right) + \frac{\gamma^{L-1}}{1-\gamma} \right) \\
 &= \beta(\delta) \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} \left( \sum_{i=1}^L \gamma^i \right) + \frac{\gamma^L}{1-\gamma} \\
 &\leq \frac{\gamma\beta(\delta)}{1-\gamma} \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \frac{\gamma^L}{1-\gamma}.
 \end{aligned}$$

The first inequality is a result of triangle inequality after adding and subtracting the term  $\gamma [P\Pi_{[0, \frac{1}{1-\gamma}]} \widehat{V}^{(L-1)}](s, a)$ . The second inequality is the error in kernel ridge regression bounded above in Term I. The third inequality bounds the difference in expectations with maximum difference. The fourth inequality is a consequence of the observation that  $V^* \in [0, \frac{1}{1-\gamma}]$  interval. The fifth inequality is obtained using the bound on the error in the proxy value function given above.

**The value function of  $\pi$  obtained from the proxy Q-function:** The following lemma establishes that the error in the value function of  $\pi$  can be bounded using the error in the proxy Q-function.

**Lemma 5.** Consider any Q-function  $\tilde{Q} : \mathcal{S} \times \mathcal{A} \rightarrow [0, \frac{1}{1-\gamma} + 1]$  satisfying  $\|\tilde{Q} - Q^*\|_\infty \leq \epsilon$ . Define policy  $\pi_{\tilde{Q}}$  such that  $\pi_{\tilde{Q}}(s) = \arg \max_{a \in \mathcal{A}} \tilde{Q}(s, a)$ . Then, for all  $s \in \mathcal{S}$ ,

$$V^*(s) - V^{\pi_{\tilde{Q}}}(s) \leq \frac{2\epsilon}{1-\gamma}.$$

Applying Lemma 5 to  $\widehat{Q}$  returned by KQLearn, it follows that with probability  $1 - \delta$ ,

$$\|V^\pi - V^*\|_\infty \leq \frac{2\gamma\beta(\delta)}{(1-\gamma)^2} \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1 + 1/\lambda^2)J}} + \frac{2\gamma^L}{(1-\gamma)^2}.$$

That completes the proof.

## B.2 Proof of Theorem 2

Theorem 2 is a consequence of Theorem 1 and the kernel specific bounds on  $\Gamma_{K,\lambda}$ . Several works have established bounds on  $\Gamma_{K,\lambda}$  for various kernels (Srinivas et al., 2010; Vakili et al., 2021c,b). We use the result in Vakili et al. (2021c) for kernels with polynomial and exponential eigendecay.

**Polynomial Eigendecay:** Consider a kernel with  $(C_p, \beta_p)$  polynomial eigendecay. We have the following bound on  $\Gamma_{K,\lambda}$  (Vakili et al., 2021c). For all  $J \in \mathbb{N}$ ,

$$\Gamma_{K,\lambda}(J) = \mathcal{O}(J^{\frac{1}{\beta_p}} \log^{1-\frac{1}{\beta_p}}(J)). \quad (28)$$

We replace this bound in the error in the value function of  $\pi$  obtained by KQLearn, in Theorem 1

$$\|V^\pi - V^*\|_\infty \leq \frac{2\gamma\beta(\delta)}{(1-\gamma)^2} \sqrt{\frac{2\Gamma_{K,\lambda}(J)}{\log(1+\frac{1}{\lambda^2})J}} + \frac{2\gamma^L}{(1-\gamma)^2}.$$

We then obtain

$$\|V^\pi - V^*\|_\infty = \mathcal{O}\left(\frac{\gamma}{(1-\gamma)^3} \left(c + \frac{1}{\lambda} \sqrt{d \log\left(\frac{Jc}{(1-\gamma)\delta}\right)}\right) \sqrt{J^{\frac{1}{\beta_p}-1} \log^{1-\frac{1}{\beta_p}}(J)}\right) + \mathcal{O}\left(\frac{\gamma^L}{(1-\gamma)^2}\right).$$

We choose  $J$  and  $L$  large enough so that each term on the right hand side is bounded by  $\epsilon/2$ . The choices of

$$\begin{aligned} J &= \mathcal{O}\left(\left(\frac{1}{\epsilon}\right)^{\frac{2\beta_p}{\beta_p-1}} \frac{\gamma^{\frac{2\beta_p}{\beta_p-1}}}{(1-\gamma)^{\frac{6\beta_p}{\beta_p-1}}} \left(c + \frac{1}{\lambda} \sqrt{d \log\left(\frac{c}{(1-\gamma)\delta}\right)}\right)^{\frac{2\beta_p}{\beta_p-1}} \log^{\frac{\beta_p}{\beta_p-1}}\left(\frac{1}{\epsilon(1-\gamma)}\right)\right), \\ L &= \mathcal{O}\left(\frac{\log(\frac{1}{\epsilon}) + \log(\frac{1}{1-\gamma})}{(1-\gamma)}\right), \end{aligned} \quad (29)$$

with proper constants ensures  $\|V^\pi - V^*\|_\infty \leq \epsilon$ . The expression can be simplified as  $J = \tilde{\mathcal{O}}\left(\frac{(\log(\frac{1}{\delta}))^{\frac{\beta_p}{\beta_p-1}}}{\epsilon^{\frac{2\beta_p}{\beta_p-1}}(1-\gamma)^{\frac{6\beta_p}{\beta_p-1}}}\right)$  and  $L = \tilde{\mathcal{O}}\left(\frac{1}{1-\gamma}\right)$ , omitting the logarithmic and constant terms. That leads to

$$N = \tilde{\mathcal{O}}\left(\frac{(\log(\frac{1}{\delta}))^{\frac{\beta_p}{\beta_p-1}}}{\epsilon^{\frac{2\beta_p}{\beta_p-1}}(1-\gamma)^{\frac{7\beta_p-1}{\beta_p-1}}}\right). \quad (30)$$

**Exponential Eigendecay:** Consider a kernel with  $(C_{e_1}, C_{e_2}, \beta_e)$  polynomial eigendecay. We have the following bound on  $\Gamma_{K,\lambda}$ . For all  $J \in \mathbb{N}$ ,

$$\Gamma_{K,\lambda}(J) = \mathcal{O}(\log^{1+\frac{1}{\beta_e}}(J)). \quad (31)$$

We replace this bound in the error in the value function of  $\pi$  obtained by KQLearn, in Theorem 1, and obtain

$$\|V^\pi - V^*\|_\infty = \mathcal{O}\left(\frac{\gamma}{(1-\gamma)^3} \left(c + \frac{1}{\lambda} \sqrt{d \log\left(\frac{Jc}{(1-\gamma)\delta}\right)}\right) \sqrt{\frac{\log^{1+\frac{1}{\beta_e}}(J)}{J}}\right) + \mathcal{O}\left(\frac{\gamma^{L-1}}{(1-\gamma)^2}\right).$$

We choose  $J$  and  $L$  large enough so that each term on the right hand side is bounded by  $\epsilon/2$ . The choices of

$$\begin{aligned} J &= \mathcal{O}\left(\left(\frac{1}{\epsilon}\right)^2 \frac{\gamma^2}{(1-\gamma)^6} \left(c + \frac{1}{\lambda} \sqrt{d \log\left(\frac{c}{(1-\gamma)\delta}\right)}\right)^2 \log^{2+\frac{1}{\beta_e}}\left(\frac{1}{\epsilon(1-\gamma)}\right)\right), \\ L &= \mathcal{O}\left(\frac{\log(\frac{1}{\epsilon}) + \log(\frac{1}{1-\gamma})}{(1-\gamma)}\right), \end{aligned} \quad (32)$$

with proper constants ensures  $\|V^\pi - V^*\|_\infty \leq \epsilon$ . The expression can be simplified as  $J = \tilde{\mathcal{O}}\left(\frac{\log(\frac{1}{\delta})}{\epsilon^2(1-\gamma)^6}\right)$  and  $L = \tilde{\mathcal{O}}\left(\frac{1}{1-\gamma}\right)$ , omitting the logarithmic and constant terms. That leads to

$$N = \tilde{\mathcal{O}}\left(\frac{\log(\frac{1}{\delta})}{\epsilon^2(1-\gamma)^7}\right). \quad (33)$$



## C PROOF OF LEMMAS

In this section, we provide the proof of auxiliary lemmas.

### C.1 Proof of Lemma 3 [RKHS norm of $PV$ ]

We have

$$\begin{aligned}
 \|PV\|_{\mathcal{H}_K} &= \left\| \int_{s' \in \mathcal{S}} P(s'|\cdot, \cdot) V(s') ds' \right\|_{\mathcal{H}_K} \\
 &\leq \int_{s' \in \mathcal{S}} \|P(s'|\cdot, \cdot) V(s')\|_{\mathcal{H}_K} ds' \\
 &= \int_{s' \in \mathcal{S}} \|P(s'|\cdot, \cdot)\|_{\mathcal{H}_K} V(s') ds' \\
 &\leq \int_{s' \in \mathcal{S}} V(s') ds'. \tag{34}
 \end{aligned}$$

where the last inequality holds by Assumption 2. We note that  $\int_{s' \in \mathcal{S}} V(s') ds' \leq \frac{1}{1-\gamma} \int_{s'} ds' \leq \frac{c}{1-\gamma}$  where  $c$  is the volume of  $\mathcal{S}$ .

### C.2 Proof of Lemma 4 [Sub-Gaussianity of $PV$ ]

The first part,  $\mathbb{E} [\widehat{PV}(s, a)] = PV(s, a)$ , follows from the definition of  $\widehat{PV}$ . For the second part note that  $\widehat{PV}$  is a random variable with a bounded support in  $[0, \frac{1}{1-\gamma}]$  by definition. Hoeffding lemma states that: let  $X$  be any random variable such that  $0 \leq X \leq a$ , then for any  $\zeta \in \mathbb{R}$ ,  $\mathbb{E} [\exp(\zeta(X - \mathbb{E}[X]))] \leq \exp(\frac{\zeta^2 a^2}{8})$ . Applying Hoeffding lemma, we can see that  $\widehat{PV}$  is sub-Gaussian with parameter  $\frac{1}{2(1-\gamma)}$ .

### C.3 Proof of Lemma 5 [Error in value function based on the error in proxy Q-function]

The proof follows similar steps as the proof of Lemma 1.11 in Agarwal et al. (2019) which considered finite state-actions. First, recall the following definitions from Section 2,

$$Q^{\pi_{\tilde{Q}}}(s, a) = r(s, a) + \gamma[PV^{\pi_{\tilde{Q}}}] (s, a) \quad \text{and} \quad Q^*(s, a) = \max_{\pi} Q(s, \pi(s)).$$

Fix state  $s \in \mathcal{S}$  and let  $\tilde{a}_s = \pi_{\tilde{Q}}(s)$ . We have

$$\begin{aligned}
 V^*(s) - V^{\pi_{\tilde{Q}}}(s) &= Q^*(s, \pi^*(s)) - Q^{\pi_{\tilde{Q}}}(s, \tilde{a}_s) \\
 &= Q^*(s, \pi^*(s)) - Q^*(s, \tilde{a}_s) + Q^*(s, \tilde{a}_s) - Q^{\pi_{\tilde{Q}}}(s, \tilde{a}_s) \\
 &= Q^*(s, \pi^*(s)) - Q^*(s, \tilde{a}_s) + \gamma[PV^*](s, \tilde{a}_s) - \gamma[PV^{\pi_{\tilde{Q}}}] (s, \tilde{a}_s) \\
 &= Q^*(s, \pi^*(s)) - Q^*(s, \tilde{a}_s) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, \tilde{a}_s)} [V^*(s') - V^{\pi_{\tilde{Q}}}(s')] \\
 &\leq Q^*(s, \pi^*(s)) - \tilde{Q}(s, \pi^*(s)) + \tilde{Q}(s, \tilde{a}_s) - Q^*(s, \tilde{a}_s) \\
 &\quad + \gamma \mathbb{E}_{s' \sim P(s, \tilde{a}_s)} [V^*(s') - V^{\pi_{\tilde{Q}}}(s')] \\
 &\leq 2 \left\| \tilde{Q} - Q^* \right\|_{\infty} + \gamma \|V^* - V^{\pi_{\tilde{Q}}}\|_{\infty},
 \end{aligned}$$

where the first inequality uses  $\tilde{Q}(s, \pi^*(s)) \leq \tilde{Q}(s, \tilde{a}_s)$  by definition of  $\tilde{a}_s$ . We thus have

$$\|V^* - V^{\pi_{\tilde{Q}}}\|_{\infty} \leq 2 \left\| \tilde{Q} - Q^* \right\|_{\infty} + \gamma \|V^* - V^{\pi_{\tilde{Q}}}\|_{\infty}. \tag{35}$$

Rearranging this inequality, we arrive at the lemma

$$\|V^* - V^{\pi_{\tilde{Q}}}\|_{\infty} \leq \frac{2}{1-\gamma} \left\| \tilde{Q} - Q^* \right\|_{\infty}. \tag{36}$$