

# Value Function Approximations via Kernel Embeddings for No-Regret Reinforcement Learning

**Sayak Ray Chowdhury**  
Microsoft Research, India

T-SAYAKR@MICROSOFT.COM

**Rafael Oliveira**  
The University of Sydney, Australia

RAFAEL.OLIVEIRA@SYDNEY.EDU.AU

**Editors:** Emtiyaz Khan and Mehmet Gönen

## Abstract

We consider the regret minimization problem in reinforcement learning (RL) in the episodic setting. In many real-world RL environments, the state and action spaces are continuous or very large. Existing approaches establish regret guarantees by either a low-dimensional representation of the stochastic transition model or an approximation of the  $Q$ -functions. However, the understanding of function approximation schemes for state-value functions largely remains missing. In this paper, we propose an online model-based RL algorithm, namely the CME-RL, that learns embeddings of the state-transition distribution in a reproducing kernel Hilbert space while carefully balancing the exploitation-exploration tradeoff. We demonstrate the efficiency of our algorithm by proving a frequentist (worst-case) regret bound that is of order  $\tilde{O}(H\gamma_N\sqrt{N})^1$ , where  $H$  is the episode length,  $N$  is the total number of time steps and  $\gamma_N$  is an information theoretic quantity relating the effective dimension of the state-action feature space. Our method bypasses the need for estimating transition probabilities and applies to any domain on which kernels can be defined. It also brings new insights into the general theory of kernel methods for approximate inference and RL regret minimization.

**Keywords:** Model-based RL; Value function approximation; Kernel mean embeddings.

## 1. Introduction

Reinforcement learning (RL) is concerned with learning to take actions to maximize rewards, by trial and error, in environments that can evolve in response to actions. A Markov decision process (MDP) (Puterman, 2014) is a popular framework to model decision making in RL environments. In the MDP, starting from an initial observed state, an agent repeatedly (a) takes an action, (b) receives a reward, and (c) observes the next state of the MDP. The traditional RL objective is a *search* goal – find a *policy* (a rule to select an action for each state) with high total reward using as few interactions with the environment as possible, also known as the sample complexity of RL (Strehl et al., 2009). This is, however, quite different from the corresponding *optimization* goal, where the learner seeks to maximize the total reward earned from all its decisions, or equivalently, minimize the *regret* or shortfall in total reward compared to that of an optimal policy (Jaksch et al., 2010). This objective is relevant in many practical sequential decision-making settings in which every decision that is taken carries utility or value – recommendation systems, sequential investment and

---

1.  $\tilde{O}(\cdot)$  hides only absolute constant and poly-logarithmic factors.

portfolio allocation, dynamic resource allocation in communication systems etc. In such *online* optimization settings, there is no separate budget or time devoted to purely exploring the unknown environment; rather, exploration and exploitation must be carefully balanced.

### 1.1. Related work

Several studies have considered the task of regret minimization in *tabular* MDPs, in which the state and action spaces are finite, and the value function is represented by a table (Jaksch et al., 2010; Osband et al., 2013; Gheshlaghi Azar et al., 2017; Dann et al., 2017; Jin et al., 2018; Efroni et al., 2019; Zanette and Brunskill, 2019). The regret bound achieved by these works essentially is proportional to  $\sqrt{SAN}$ , where  $S$  and  $A$  denote the numbers of states and actions, respectively, and  $N$  the total number of steps. In many practical applications, however, the number of states and actions is enormous. For example, the game of Go has a state space with size  $3^{361}$ , and the state and action spaces of certain robotics applications can even be continuous. These continuous state and action spaces make RL a challenging task, especially in terms of generalizing learnt knowledge across unseen states and actions. In such cases, the tabular model suffers from the “curse of dimensionality” problem. To tackle this issue, the popular “optimism in the face of uncertainty” principle from Jaksch et al. (2010) has been extended to handle continuous MDPs, when assuming some Lipschitz-like smoothness or regularity on the rewards and dynamics (Ortner and Ryabko, 2012; Domingues et al., 2020).

Another line of work considers *function approximation*, i.e., they use features to parameterize reward and transition models, with the hope that the features can capture leading structures of the MDP (Osband and Van Roy, 2014; Chowdhury and Gopalan, 2019). The *model-based* algorithms developed in these works assume oracle access to an optimistic planner to facilitate the learning. The optimistic planning step is quite prohibitive and often becomes computationally intractable for continuous state and action spaces. Yang and Wang (2019) consider a low-rank bilinear transition model bypassing the complicated planning step; however, their algorithm potentially needs to compute the *value function* across all states. This suffers an  $\Omega(S)$  computational complexity and as a consequence cannot directly handle continuous state spaces. Ayoub et al. (2020) consider linear-mixture transition structure that includes the bilinear model as a special case. However, their algorithm too suffers the  $\Omega(S)$  computational complexity. To alleviate the computational burden intrinsic to these model-based approaches, a recent body of work parameterizes the *value functions* directly, using  $d$ -dimensional state-action feature maps, and develop *model-free* algorithms bypassing the need for fully learning the reward and transition models (Jin et al., 2019; Wang et al., 2019; Zanette et al., 2020a). Under the assumption that the (action-)value function can be approximated by a linear or a generalized linear function of the feature vectors, these papers develop algorithms with regret bound proportional to  $\text{poly}(d)\sqrt{T}$ , which is independent of the size of the state and action spaces. Wang et al. (2020) generalizes this approach by designing an algorithm that works with general (non-linear) value function approximators and prove a similar regret guarantee that depends on the eluder dimension (Russo and Van Roy, 2013) and log-covering number of the underlying function class.

A few recent works have proposed kernel-based value function approximation algorithms. Yang et al. (2020) consider kernel and neural function approximations and designed algo-

gorithms with regret characterized by intrinsic complexity of the function classes. More closely related to our work, Domingues et al. (2021) recently proposed a kernel-based RL algorithm via value function approximation. Their main assumption relies on Lipschitz continuity of the reward functions and the state transition kernels. In contrast to their work, we are able to obtain tighter regret bounds by applying typical assumptions in the kernel embeddings literature, which we show are satisfied for a variety of practical systems. Nevertheless, there is a lack of theoretical understanding in designing provably efficient model-based RL algorithms with (non-linear) value function approximation, which we aim to address.

## 1.2. Contributions

In this work, we revisit function approximation in RL by modeling the value functions as elements of a reproducing kernel Hilbert space (RKHS) (Schölkopf and Smola, 2002) compatible with a (possibly infinite dimensional) state feature map. The main motivation behind this formulation is that the conditional expectations of any function in the RKHS become a linear operation, via the RKHS inner product with an appropriate distribution embedding, known as the *conditional mean embedding* (Muandet et al., 2016). In recent years, conditional mean embeddings (CMEs) have found extensive applications in many machine learning tasks (Song et al., 2009, 2010a,b, 2013; Fukumizu et al., 2008, 2009; Hsu and Ramos, 2019; Chowdhury et al., 2020). The foremost advantage of CMEs in our setup is that one can directly compute conditional expectations of the value functions based only on the observed data, since the alternative approach of estimating the transition probabilities as an intermediate step scales poorly with the dimension of the state space (Grünwälder et al., 2012). The convergence of conditional mean estimates to the true embeddings in the RKHS norm has been established by Grünwälder et al. (2012) assuming access to *independent and identically distributed* (i.i.d.) transition samples (the “simulator” setting). However, in the online RL environment like the one considered in this work, one collects data based on past observations, and hence the existing result fails to remain useful. Against this backdrop, we make the following contributions:

- In online RL environment, we derive a concentration inequality for mean embedding estimates of the transition distribution around the true embeddings as a function of the uncertainties around these estimates (Theorem 1). This bound serves as a key tool in designing our model-based RL algorithm, while also being of independent interest.
- Focusing on the value function approximation in the RKHS setting, we present the first model-based RL algorithm, namely the *Conditional Mean Embedding RL* (CME-RL), that is provably efficient in regret performance and does not require any additional oracle access or stronger computational assumptions (Algorithm 1). Concretely, in the general episodic MDP setting, CME-RL enjoys a regret bound of  $\tilde{O}(H\gamma_N\sqrt{N})$ , where  $H$  is the length of each episode,  $\gamma_N$  is a complexity measure relating the effective dimension of the RKHS compatible with the state-action features (Theorem 2).
- Our approach is also robust to the RKHS modelling assumption: when the value functions are not elements of the RKHS, but  $\zeta$ -close to some RKHS element in the  $\ell_\infty$  norm, then (a modified version of) CME-RL achieves a  $\tilde{O}(H\gamma_N\sqrt{N}+\zeta N)$  regret, where the linear regret term arises due to the function class misspecification (Theorem 3).

## 2. Preliminaries

**Notations** We begin by introducing some notations. Let  $\mathcal{H}$  be an arbitrary Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  and corresponding norm  $\|\cdot\|_{\mathcal{H}}$ . When  $\mathcal{G}$  is another Hilbert space, we denote by  $\mathcal{L}(\mathcal{G}, \mathcal{H})$  the Banach space of linear operators  $A : \mathcal{G} \rightarrow \mathcal{H}$  with bounded operator norm  $\|A\| := \sup_{\|g\|_{\mathcal{G}}=1} \|A g\|_{\mathcal{H}}$ . We let  $\text{HS}(\mathcal{G}, \mathcal{H})$  denote the subspace of operators in  $\mathcal{L}(\mathcal{G}, \mathcal{H})$  with bounded Hilbert-Schmidt norm, defined for  $A \in \text{HS}(\mathcal{G}, \mathcal{H})$  as  $\|A\|_{\text{HS}} := \left( \sum_{i,j=1}^{\infty} \langle f_i, A g_j \rangle_{\mathcal{H}}^2 \right)^{1/2}$ , where the  $f_i$ 's form a complete orthonormal system (CONS) for  $\mathcal{H}$  and the  $g_j$ 's form a CONS for  $\mathcal{G}$ . In the case  $\mathcal{G} = \mathcal{H}$ , we set  $\mathcal{L}(\mathcal{H}) := \mathcal{L}(\mathcal{H}, \mathcal{H})$ . We denote by  $\mathcal{L}_+(\mathcal{H})$  the set of all bounded, positive-definite linear operators on  $\mathcal{H}$ , i.e.,  $A \in \mathcal{L}_+(\mathcal{H})$  if, for any non-zero  $h \in \mathcal{H}$ ,  $\langle h, A h \rangle_{\mathcal{H}} > 0$ .

**Regret minimization in finite-horizon episodic MDPs** We consider episodic reinforcement learning in a finite-horizon Markov decision process (MDP) of episode length  $H$  with (possibly infinite) state and action spaces  $\mathcal{S}$  and  $\mathcal{A}$ , respectively, reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ , and transition probability measure  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , where  $\Delta(\mathcal{S})$  denotes the probability simplex on  $\mathcal{S}$ . The learning agent interacts with the MDP in episodes and, at each episode  $t$ , a trajectory  $(s_1^t, a_1^t, r_1^t, \dots, s_H^t, a_H^t, r_H^t, s_{H+1}^t)$  is generated. Here  $a_h^t$  denotes the action taken at state  $s_h^t$ ,  $r_h^t := R(s_h^t, a_h^t)$  denotes the immediate reward, and  $s_{h+1}^t \sim P(\cdot | s_h^t, a_h^t)$  denotes the random next state. The initial state  $s_1^t$  is assumed to be fixed and history independent, and can even be possibly chosen by an adversary. The episode terminates when  $s_{H+1}^t$  is reached, where the agent cannot take any action and hence receives no reward. The actions are chosen following some policy  $\pi = (\pi_1, \dots, \pi_H)$ , where each  $\pi_h$  is a mapping from the state space  $\mathcal{S}$  into the action space  $\mathcal{A}$ . The agent would like to find a policy  $\pi$  that maximizes the long-term expected cumulative reward starting from every state  $s \in \mathcal{S}$  and every step  $h \in [H]$ , defined as:

$$V_h^\pi(s) := \mathbb{E} \left[ \sum_{j=h}^H R(s_j, \pi_j(s_j)) \mid s_h = s \right].$$

We call  $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$  the value function of policy  $\pi$  at step  $h$ . Accordingly, we also define the action-value function, or  $Q$ -function,  $Q_h^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  as:

$$Q_h^\pi(s, a) := R(s, a) + \mathbb{E} \left[ \sum_{j=h+1}^H R(s_j, \pi_j(s_j)) \mid s_h = s, a_h = a \right],$$

which gives the expected value of cumulative rewards starting from a state-action pair at the  $h$ -th step and following the policy  $\pi$  afterwards. Note that  $V_h^\pi(s) = Q_h^\pi(s, \pi_h(s))$  and it satisfies the Bellman equation:

$$V_h^\pi(s) = R(s, \pi_h(s)) + \mathbb{E}_{X \sim P(\cdot | s, \pi_h(s))} [V_{h+1}^\pi(X)], \quad \forall h \in [H], \quad (1)$$

with  $V_{H+1}^\pi(s) = 0$  for all  $s \in \mathcal{S}$ . We denote by  $\pi^*$  an optimal policy satisfying:

$$V_h^{\pi^*}(s) = \max_{\pi \in \Pi} V_h^\pi(s), \quad \forall s \in \mathcal{S}, \forall h \in [H],$$

where  $\Pi$  is the set of all non-stationary policies. Since the episode length is finite, such a policy exists when the action space  $\mathcal{A}$  is large but finite (Puterman, 2014). We denote the optimal value function by  $V_h^*(s) := V_h^{\pi^*}(s)$ . We also denote the optimal action-value function (or  $Q$ -function) as  $Q_h^*(s, a) = \max_{\pi} Q_h^\pi(s, a)$ . It is easily shown that the optimal

action-value function satisfies the Bellman optimality equation:

$$Q_h^*(s, a) := R(s, a) + \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^*(X)], \quad \forall h \in [H], \quad (2)$$

with  $V_h^*(s) = \max_{a \in \mathcal{A}} Q_h^*(s, a)$ . This implies that the optimal policy is the greedy policy with respect to the optimal action-value functions. Thus, to find the optimal policy  $\pi^*$ , it suffices to estimate the optimal action-value functions  $(Q_h^*)_{h \in [H]}$ .

The agent aims to learn the optimal policy by interacting with the environment during a set of episodes. We measure performance of the agent by the cumulative (pseudo) regret accumulated over  $T$  episodes, defined as:

$$\mathcal{R}(N) := \sum_{t=1}^T \left[ V_1^*(s_1^t) - V_1^{\pi^t}(s_1^t) \right],$$

where  $\pi^t$  is the policy chosen by the agent at episode  $t$  and  $N = TH$  is the total number of steps. The regret measures the quantum of reward that the learner gives up by not knowing the MDP in advance and applying the optimal policy  $\pi^*$  from the start. We seek algorithms that attain sublinear regret  $\mathcal{R}(N) = o(N)$  in the number of steps they face, since, for instance, an algorithm that does not adapt its policy selection behavior depending on past experience can easily be seen to achieve linear ( $\Omega(N)$ ) regret (Lai and Robbins, 1985).

**Value function approximation in episodic MDPs** A very large or possibly infinite state and action space makes reinforcement learning a challenging task. To obtain sub-linear regret guarantees, it is necessary to posit some regularity assumptions on the underlying function class. In this paper, we use reproducing kernel Hilbert spaces to model the value functions. Let  $\mathcal{H}_\psi$  and  $\mathcal{H}_\varphi$  be two RKHSs with continuous positive semi-definite kernel functions  $k_\psi : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}_+$  and  $k_\varphi : (\mathcal{S} \times \mathcal{A}) \times (\mathcal{S} \times \mathcal{A}) \rightarrow \mathbb{R}_+$ , with corresponding inner products  $\langle \cdot, \cdot \rangle_{\mathcal{H}_\psi}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{H}_\varphi}$ , respectively. There exist feature maps  $\psi : \mathcal{S} \rightarrow \mathcal{H}_\psi$  and  $\varphi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{H}_\varphi$  such that  $k_\psi(\cdot, \cdot) = \langle \psi(\cdot), \psi(\cdot) \rangle_{\mathcal{H}_\psi}$  and  $k_\varphi(\cdot, \cdot) = \langle \varphi(\cdot), \varphi(\cdot) \rangle_{\mathcal{H}_\varphi}$ , respectively (Steinwart and Christmann, 2008).

The weakest assumption one can pose on the value functions is realizability, which posits that the optimal value functions  $(V_h^*)_{h \in [H]}$  lie in the RKHS  $\mathcal{H}_\psi$ , or at least are well-approximated by  $\mathcal{H}_\psi$ . For stateless MDPs or multi-armed bandits where  $H = 1$ , realizability alone suffices for provably efficient algorithms (Abbasi-Yadkori et al., 2011; Chowdhury and Gopalan, 2017). But it does not seem to be sufficient when  $H > 1$ , and in these settings it is common to make stronger assumptions (Jin et al., 2019; Wang et al., 2019, 2020). Following these works, our main assumption is a closure property for all value functions in the following class:

$$\mathcal{V} := \left\{ s \mapsto \min \left\{ H, \max_{a \in \mathcal{A}} \left\{ R(s, a) + \langle \varphi(s, a), \mu \rangle_{\mathcal{H}_\varphi} + \eta \sqrt{\langle \varphi(s, a), \Sigma^{-1} \varphi(s, a) \rangle_{\mathcal{H}_\varphi}} \right\} \right\} \right\}, \quad (3)$$

where  $0 < \eta < \infty$ ,  $\mu \in \mathcal{H}_\varphi$  and  $\Sigma \in \mathcal{L}_+(\mathcal{H}_\varphi)$  are the parameters of the function class.

**Assumption 1 (Optimistic closure)** For any  $V \in \mathcal{V}$  (cf. Equation 3), we have  $V \in \mathcal{H}_\psi$ . Furthermore, for a positive constant  $B_V$ , we have  $\|V\|_{\mathcal{H}_\psi} \leq B_V$ .

While this property seems quite strong, we note that related closure-type assumptions are common in the literature. We will relax this assumption later in Section 4.3. In addition, our results do not require explicit knowledge of  $\mathcal{H}_\psi$  nor its kernel  $k_\psi$ , as we will only interact with elements of  $\mathcal{V}$  via point evaluations and RKHS norm bounds.

### 3. RKHS embeddings of transition distribution

In order to find an estimate of the optimal value function, it is imperative to estimate the conditional expectations of the form  $\mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)]$ . In the model-based approach considered in this work, we do so by estimating the mean embedding of the conditional distribution  $P(\cdot|s,a)$ , which is the focus of this section. For a bounded kernel<sup>2</sup>  $k_\psi$  on the state space  $\mathcal{S}$ , the mean embedding of the conditional distribution  $P(\cdot|s,a)$  in  $\mathcal{H}_\psi$  is an element  $\vartheta_P^{(s,a)} \in \mathcal{H}_\psi$  such that:

$$\forall f \in \mathcal{H}_\psi, \quad \mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)] = \langle f, \vartheta_P^{(s,a)} \rangle_{\mathcal{H}_\psi}. \quad (4)$$

The mean embedding can be explicitly expressed as a function:

$$\vartheta_P^{(s,a)}(y) = \mathbb{E}_{X \sim P(\cdot|s,a)}[k_\psi(X, y)],$$

for all  $y \in \mathcal{S}$ . If the kernel  $k_\psi$  is characteristic, such as a stationary kernel, then the mapping  $P(\cdot|s,a) \mapsto \vartheta_P^{(s,a)}$  is injective, defining a one-to-one relationship between transition distributions and elements of  $\mathcal{H}_\psi$  (Sriperumbudur et al., 2011). Following existing works (Song et al., 2009; Grünewälder et al., 2012), we now make a smoothness assumption on the transition distribution.

**Assumption 2** For any  $f \in \mathcal{H}_\psi$ , the function  $(s, a) \mapsto \mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)]$  lies in  $\mathcal{H}_\varphi$ .

Under Assumption 2, the mean embeddings admit a linear representation in state-action features via the conditional embedding operator  $\Theta_P \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\psi)$  such that:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad \vartheta_P^{(s,a)} = \Theta_P \varphi(s, a). \quad (5)$$

Assumption 2 always holds for finite domains with characteristic kernels. Though it is not necessarily true for continuous domains, we note that the CMEs for classical linear (Abbasi-Yadkori and Szepesvári, 2011) and non-linear (Kakade et al., 2020) dynamical systems satisfy this assumption.

#### 3.1. Sample estimate of conditional mean embedding

At the beginning of each episode  $t$ , given the observations  $\mathcal{D}_t := (s_h^\tau, a_h^\tau, s_{h+1}^\tau)_{\tau < t, h \leq H}$  until episode  $t - 1$ , we consider a sample based estimate of the conditional embedding operator. This is achieved by solving the following ridge-regression problem:

$$\min_{\Theta \in \text{HS}(\mathcal{H}_\varphi, \mathcal{H}_\psi)} \sum_{\tau < t, h \leq H} \|\psi(s_{h+1}^\tau) - \Theta \varphi(s_h^\tau, a_h^\tau)\|_{\mathcal{H}_\psi}^2 + \lambda \|\Theta\|_{\text{HS}}^2, \quad (6)$$

where  $\lambda > 0$  is a regularising constant. The solution of Equation 6 is given by:

$$\hat{\Theta}_t = \sum_{\tau < t, h \leq H} \psi(s_{h+1}^\tau) \otimes \varphi(s_h^\tau, a_h^\tau) \left( \hat{\mathbb{C}}_{\varphi,t} + \lambda \mathbf{I} \right)^{-1}, \quad (7)$$

where  $\hat{\mathbb{C}}_{\varphi,t} := \sum_{\tau < t, h \leq H} \varphi(s_h^\tau, a_h^\tau) \otimes \varphi(s_h^\tau, a_h^\tau)$  and  $\otimes$  denotes the tensor product of elements in a Hilbert space. To simplify notations, we now let  $n = (t - 1)H$  denote the total number of steps completed at the beginning of episode  $t$ . We denote a vector  $k_{\varphi,t}(s, a) \in \mathbb{R}^n$  and a matrix  $\mathbb{K}_{\varphi,t} \in \mathbb{R}^{n \times n}$  by:

$$k_{\varphi,t}(s, a) := [k_\varphi((s_h^\tau, a_h^\tau), (s, a))]_{\tau < t, h \leq H}, \quad \mathbb{K}_{\varphi,t} := [k_\varphi((s_h^\tau, a_h^\tau), (s_{h'}^{\tau'}, a_{h'}^{\tau'}))]_{\tau, \tau' < t, h, h' \leq H}.$$

2. Boundedness of a kernel holds for any stationary kernel, e.g., the squared exponential kernel and the Matérn kernel (Rasmussen and Williams, 2006).



Then, via [Equation 7](#), the conditional mean embeddings can be estimated as

$$\widehat{\vartheta}_t^{(s,a)} = \widehat{\Theta}_t \varphi(s, a) = \sum_{\tau < t, h \leq H} [\alpha_t(s, a)]_{(\tau, h)} \psi(s_{h+1}^\tau), \quad (8)$$

where we define the weight vector  $\alpha_t(s, a) := (\mathbf{K}_{\varphi, t} + \lambda \mathbf{I})^{-1} k_{\varphi, t}(s, a)$ .

### 3.2. Concentration of mean embedding estimates

In this section, we show that for any state-action pair  $(s, a)$ , the CME estimates  $\widehat{\vartheta}_t^{(s,a)}$  lies within a high-probability confidence region around the true embedding  $\vartheta_P^{(s,a)}$ . This eventually translates, via [Equation 4](#), to a concentration property of  $\langle f, \widehat{\vartheta}_t^{(s,a)} \rangle_{\mathcal{H}_\psi}$  around  $\mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)]$  for any  $f \in \mathcal{H}_\psi$ . The uncertainty of CME estimates can be characterized by the variance estimate  $\sigma_{\varphi, t}^2(s, a) := \lambda \langle \varphi(s, a), \mathbf{M}_t^{-1} \varphi(s, a) \rangle_{\mathcal{H}_\varphi}$ , where  $\mathbf{M}_t := \widehat{\mathbf{C}}_{\varphi, t} + \lambda \mathbf{I}$ . To see this, note that an application of Sherman-Morrison formula yields:

$$\sigma_{\varphi, t}^2(s, a) := k_\varphi((s, a), (s, a)) - k_{\varphi, t}(s, a)^\top (\mathbf{K}_{\varphi, t} + \lambda \mathbf{I})^{-1} k_{\varphi, t}(s, a), \quad (9)$$

which is equivalent to the predictive variance of a Gaussian process (GP) ([Rasmussen and Williams, 2006](#)). Although a sample from a GP is usually not an element of the RKHS defined by its kernel ([Lukic and Beder, 2001](#)), the following result allows us to use  $\sigma_{\varphi, t}^2(s, a)$  as an error measure.

**Theorem 1 (Concentration of the conditional embedding operator)** *Suppose that  $\sup_{s \in \mathcal{S}} \sqrt{k_\psi(s, s)} \leq B_\psi$ . Then, under [Assumption 2](#), for any  $\lambda > 0$  and  $\delta \in (0, 1]$ ,*

$$\mathbb{P} \left[ \forall t \in \mathbb{N}, \left\| (\Theta_P - \widehat{\Theta}_t) \mathbf{M}_t^{1/2} \right\| \leq \beta_t(\delta) \right] \geq 1 - \delta,$$

where  $\beta_t(\delta) := \sqrt{2\lambda B_P^2 + 256(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi, t})^{1/2}) \log(2t^2 H / \delta)}$ ,  $B_P \geq \|\Theta_P\|_{\text{HS}}$ .

[Theorem 1](#) implies a concentration inequality for the CME estimates, since, for all  $t \geq 1$ :

$$\left\| \vartheta_P^{(s,a)} - \widehat{\vartheta}_t^{(s,a)} \right\|_{\mathcal{H}_\psi} \leq \left\| (\Theta_P - \widehat{\Theta}_t) \mathbf{M}_t^{1/2} \right\| \|\varphi(s, a)\|_{\mathbf{M}_t^{-1}} \leq \beta_t(\delta) \lambda^{-1/2} \sigma_{\varphi, t}(s, a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A},$$

with probability at least  $1 - \delta$ . This forms the core of our value function approximations.<sup>3</sup>

**Remark 1** *Considering the simulation setting, [Grünwälder et al. \(2012\)](#) assume access to a sample  $(s_i, a_i, s'_i)_{i=1}^m$ , drawn i.i.d. from a joint distribution  $P_0$  such that the conditional probabilities satisfy  $P_0(s'_i | s_i, a_i) = P(s'_i | s_i, a_i), \forall i$ . Under [Assumption 2](#), they establish the convergence of CME estimates  $\widehat{\vartheta}_t^{(s,a)}$  to the true CMEs  $\vartheta_P^{(s,a)}$  in  $P_0$ -probability. This guarantee, however, does not apply to our setting, since we do not assume any simulator access.*

**Proof sketch of [Theorem 1](#)** To derive this result, we note that the sequence evaluation noise  $\varepsilon_h^t := \psi(s_{h+1}^t) - \Theta_P \varphi(s_h^t, a_h^t)$  at each step  $h$  of episode  $t$  forms a martingale difference sequence, with each element having a bounded RKHS norm. We overload notation to define, for each pair  $(t, h)$ , the operator  $\mathbf{M}_{t, h} = \mathbf{M}_t + \sum_{j \leq h} \varphi(s_j^t, a_j^t) \otimes \varphi(s_j^\tau, a_j^\tau)$ , and the estimate

$$\widehat{\Theta}_{t, h} = \left( \sum_{\tau < t, h \leq H} \psi(s_{h+1}^\tau) \otimes \varphi(s_h^\tau, a_h^\tau) + \sum_{j \leq h} \psi(s_{j+1}^t) \otimes \varphi(s_j^t, a_j^t) \right) \mathbf{M}_{t, h}^{-1}.$$

3. [Deshmukh et al. \(2017\)](#) employ a variant of kernel ridge regression to learn the mean reward function of a contextual bandit problem. Their concentration bound holds only for finite action space. In contrast, [Theorem 1](#) holds for infinite state-action spaces, and hence, can be seen as a generalization of their result.

Now, we consider the random variable  $z_{t,h} = \left\| (\widehat{\Theta}_{t,h} - \Theta_P) M_{t,h}^{1/2} \right\|_{\text{HS}}^2$ , and prove a high-probability upper bound on it using Azuma-Hoeffding's inequality for martingales. In fact, we show that  $z_{t,h} \leq \beta_{t,h}^2(\delta)$  uniformly over all pair  $(t, h)$  with probability at least  $1 - \delta$ , where  $\beta_{t,h}(\delta)$  is defined similarly to  $\beta_t(\delta)$  with only  $K_{\varphi,t}$  being replaced by  $K_{\varphi,t,h} := [k_{\varphi}((s_j^{\tau}, a_j^{\tau}), (s_{j'}^{\tau}, a_{j'}^{\tau}))]_{(\tau,j),(\tau',j') \leq (t,h)}$  – the gram-matrix at step  $h$  of episode  $t$ . The proof then follows by noting that  $\left\| (\widehat{\Theta}_t - \Theta_P) M_t^{1/2} \right\|_{\text{HS}} = z_{t-1,H}^{1/2} \leq \beta_{t-1,H}(\delta) \leq \beta_t(\delta)$ . The complete proof is given in the supplementary material.

## 4. RL exploration using RKHS embeddings

In this section, we aim to develop an online RL algorithm using the conditional mean embedding estimates that balances exploration and exploitation (near) optimally. We realize this, at a high level, by following the Upper-Confidence Bound (UCB) principle and thus our algorithm falls in a similar framework as in [Jaksch et al. \(2010\)](#); [Gheshlaghi Azar et al. \(2017\)](#); [Yang and Wang \(2019\)](#).

### 4.1. The Conditional Mean Embedding RL (CME-RL) algorithm

At a high level, each episode  $t$  consists of two passes over all steps. In the first pass, we maintain the  $Q$ -function estimates via dynamic programming. To balance the exploration-exploitation trade-off, we first define a confidence set  $\mathcal{C}_t$  that contains the set of conditional embedding operators that are deemed to be consistent with all the data that has been collected in the past. Specifically, for any  $\delta \in (0, 1]$ ,  $\lambda > 0$  and constants  $B_P$  and  $B_{\psi}$ , [Theorem 1](#) governs us to define the confidence set

$$\mathcal{C}_t := \left\{ \Theta \in \mathcal{L}(\mathcal{H}_{\varphi}, \mathcal{H}_{\psi}) : \left\| (\Theta - \widehat{\Theta}_t) M_t^{1/2} \right\| \leq \beta_t(\delta/2) \right\}, \quad (10)$$

where  $\beta_t(\cdot)$  governs the exploration-exploitation trade-off. This confidence set is then used to compute the optimistic  $Q$ -estimates, starting with  $V_{H+1}^t(s) = 0$ , and setting:

$$\text{for } h = H, H-1, \dots, 1, \quad V_h^t(s) = \min \left\{ H, \max_{a \in \mathcal{A}} Q_h^t(s, a) \right\}, \quad (11)$$

$$Q_h^t(s, a) = R(s, a) + \max_{\Theta_{P'} \in \mathcal{C}_t} \mathbb{E}_{X \sim P'(\cdot|s,a)} [V_{h+1}^t(X)]. \quad (12)$$

We note here that we only require an optimistic estimate of the optimal  $Q$ -function. Hence, it is not necessary to solve the maximization problem in [Equation 12](#) explicitly. In fact, we can use a closed-form expression instead of searching for the optimal embedding operator  $\Theta_{P'}$  in the confidence set  $\mathcal{C}_t$ . If the value estimate  $V_{h+1}^t$  lies in the RKHS  $\mathcal{H}_{\psi}$ , we then have from [Equation 4](#) that  $\mathbb{E}_{X \sim P'(\cdot|s,a)} [V_{h+1}^t(X)] = \langle V_{h+1}^t, \vartheta_{P'}^{s,a} \rangle_{\mathcal{H}_{\psi}}$ , and from [Equation 8](#) that:

$$\left\langle V_{h+1}^t, \widehat{\vartheta}_t^{(s,a)} \right\rangle_{\mathcal{H}_{\psi}} = \alpha_t(s, a)^{\top} v_{h+1}^t = k_{\varphi,t}(s, a)^{\top} (K_{\varphi,t} + \lambda \mathbf{I})^{-1} v_{h+1}^t,$$

where we define the vector  $v_{h+1}^t := [V_{h+1}^t(s_{h'+1}^{\tau})]_{\tau < t, h' \leq H}$ . Now, since the confidence set  $\mathcal{C}_t$  is convex, the  $Q$ -updates given by [Equation 12](#) admit the closed-form expression:

$$Q_h^t(s, a) = R(s, a) + k_{\varphi,t}(s, a)^{\top} (K_{\varphi,t} + \lambda \mathbf{I})^{-1} v_{h+1}^t + \|V_{h+1}^t\|_{\mathcal{H}_{\psi}} \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi,t}(s, a). \quad (13)$$

We now note that, by the optimistic closure property ([Assumption 1](#)), the value estimate  $V_h^t$  given by [Equation 11](#) lies in the RKHS  $\mathcal{H}_{\psi}$ , rendering the closed-form expression in [Equation 13](#) valid.



In the second pass, we execute the greedy policy with respect to the  $Q$ -function estimates obtained in the first pass. Specifically, at each step  $h$ , we chose the action:

$$a_h^t = \pi_h^t(s_h^t) \in \operatorname{argmax}_{a \in \mathcal{A}} Q_h^t(s_h^t, a) . \quad (14)$$

The pseudo-code of CME-RL is given in [Algorithm 1](#). Note that, in order to implement CME-RL, we do not need to know the kernel  $k_\psi$ ; only the knowledge of the upper bound  $B_V$  over the RKHS norm of  $V_{h+1}^t$  suffices our purpose. For simplicity of representation, we assume that the agent, while not knowing the conditional mean embedding operator  $\Theta_P$ , knows the reward function  $R$ . When  $R$  is unknown but an element of the RKHS  $\mathcal{H}_\varphi$ , our algorithm can be extended naturally with an optimistic reward estimation step at each episode, similar to the contextual bandit setting ([Chowdhury and Gopalan, 2017](#)).

---

**Algorithm 1:** Conditional Mean Embedding RL (CME-RL)

---

```

1 Input: Kernel  $k_\varphi$ , constants  $B_P$ ,  $B_V$  and  $B_\psi$ , parameters  $\eta > 0$  and  $\delta \in (0, 1]$ 
2 for episode  $t = 1, \dots, T$  do
3   Receive the initial state  $s_1^t$  and set  $V_{H+1}^t(\cdot) = 0$ 
4   for step  $h = H, \dots, 1$  do // Update value function estimates
5      $Q_h^t(\cdot, \cdot) = R(\cdot, \cdot) + k_{\varphi,t}(\cdot, \cdot)^\top (\mathbf{K}_{\varphi,t} + \lambda \mathbf{I})^{-1} v_{h+1}^t + B_V \beta_t (\delta/2) \lambda^{-1/2} \sigma_{\varphi,t}(\cdot, \cdot)$ 
6      $V_h^t(\cdot) = \min \{H, \max_{a \in \mathcal{A}} Q_h^t(\cdot, a)\}$ 
7   for step  $h = 1, \dots, H$  do // Run episode
8     Take action  $a_h^t \in \operatorname{argmax}_{a \in \mathcal{A}} Q_h^t(s_h^t, a)$  and observe next state  $s_{h+1}^t \sim P(\cdot | s_h^t, a_h^t)$ 

```

---

**Computational complexity of CME-RL** The dominant cost is evaluating the  $Q$ -function estimates  $Q_h^t$  ([Equation 13](#)). As typical in kernel methods ([Schölkopf and Smola, 2002](#)), it involves inversion of  $tH \times tH$  matrices, which take  $O(t^3 H^3)$  time. In the policy execution phase ([Equation 14](#)), we do not need to compute the entire  $Q$ -function as the algorithm only queries  $Q$ -values at visited states. Hence, assuming a constant cost of optimizing over the actions, the per-episode running time is  $O(t^3 H^4)$ . However, using standard sketching techniques like the Nyström approximation ([Drineas and Mahoney, 2005](#)) or the random Fourier features approximation ([Rahimi and Recht, 2007](#)), and by using the Sherman-Morrison formula to amortize matrix inversions, per-episode running cost can be reduced to  $O(m^2 H)$ , where  $m$  is the dimension of feature approximations.

## 4.2. Regret bound for CME-RL

In this section, we present the regret guarantee of our algorithm. We first define

$$\gamma_N \equiv \gamma_{\varphi,\lambda,N} := \sup_{\mathcal{X} \subset \mathcal{S} \times \mathcal{A}: |\mathcal{X}|=N} \frac{1}{2} \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,\mathcal{X}}) ,$$

where  $\mathcal{X} = \{(s_i, a_i)\}_{i \in [N]}$  and  $\mathbf{K}_{\varphi,\mathcal{X}} = [k_\varphi((s_i, a_i), (s_j, a_j))]_{i,j \in [N]}$  is the gram matrix over the data set  $\mathcal{X}$ .  $\gamma_N$  denotes the *maximum information gain* about a (random) function  $f$  sampled from a zero-mean GP with covariance function  $k_\varphi$  after  $N$  noisy observations, obtained by passing  $f$  through an i.i.d. Gaussian channel  $\mathcal{N}(0, \lambda)$ . Consider the case when  $k_\varphi$  is a squared exponential kernel on  $\mathbb{R}^d$ . Then it can be verified that  $\gamma_N = O((\log N)^{d+1})$  ([Srinivas et al., 2009](#)).

**Theorem 2 (Cumulative regret of CME-RL)** *Under assumptions 1 and 2, after interacting with the environment for  $N = TH$  steps, with probability at least  $1 - \delta$ , CME-RL (Algorithm 1) achieves the regret bound*

$$\mathcal{R}(N) \leq 2B_V \alpha_{N,\delta} \sqrt{2(1 + \lambda^{-1} B_\varphi^2 H) N \gamma_N} + 2H \sqrt{2N \log(2/\delta)},$$

where  $B_\varphi \geq \sup_{s,a} \sqrt{k_\varphi((s,a), (s,a))}$ , and  $\alpha_{N,\delta} := \sqrt{2\lambda B_P^2 + 256(1 + \lambda^{-1})\gamma_N \log(4N^2/\delta)}$ .

Theorem 2 yields a  $\tilde{O}(H\gamma_N\sqrt{N})$  regret bound for CME-RL. Comparing to the minimax regret in tabular setting,  $\Theta(H\sqrt{SAN})$  (Gheshlaghi Azar et al., 2017), our bound replaces the sublinear dependency on the number of state-action pairs by a linear dependency on the intrinsic complexity measure,  $\gamma_N$ , of the feature space  $\mathcal{H}_\varphi$ , which is crucial in the large state-action space setting that entails function approximation. Additionally, in the kernelized bandit setting ( $H = 1$ ), our bound matches the best known upper bound  $O(\gamma_N\sqrt{N})$  (Chowdhury and Gopalan, 2017). We note, however, that while an MDP has state transitions, the bandits do not, and a naive adaptation of existing kernelized bandit algorithms to this setting would give a regret exponential in episode length  $H$ . Furthermore, due to the Markov transition structure, the lower bound for kernelized bandits (Scarlett et al., 2017) does not directly apply here. Hence, it remains an interesting future direction to determine the optimal dependency on  $\gamma_N$ .

**Conversion to PAC guarantee** Similarly to the discussion in Jin et al. (2019), our regret bound directly translates to a sample complexity or probably approximately correct (PAC) guarantee in the following sense. Assuming a fixed initial state  $s_1^t = s$  for each episode  $t$ , with at least a constant probability, we can learn an  $\varepsilon$ -optimal policy  $\pi$  that satisfies  $V_1^*(s) - V_1^\pi(s) \leq \varepsilon$  by running CME-RL for  $T = O(d_{\text{eff}}^2 H^2 / \varepsilon^2)$  episodes, where  $d_{\text{eff}}$  is a known upper bound over  $\gamma_N$ , and then output the greedy policy according to the  $Q$ -function at  $t$ -th episode, where  $t$  is sampled uniformly from  $[T]$ . Here  $d_{\text{eff}}$  effectively captures the number of significant dimensions of  $\mathcal{H}_\varphi$ .

**Remark 2 Yang and Wang (2019)** *assumes the model  $P(s'|s, a) = \langle \psi(s'), \Theta_P \varphi(s, a) \rangle_{\mathcal{H}_\psi}$ , and propose an algorithm with regret  $\tilde{O}(H^2 \gamma_N \sqrt{N})$ . In comparison, we get an  $O(H)$  factor improvement thanks to a tighter control over the sum of predictive variances. Furthermore, their algorithm can't be implemented exactly as they need to apply random sampling to approximate the estimate  $\hat{\Theta}_t$ . We overcome this implementational bottleneck by virtue of our novel confidence set construction using the CME estimates (Theorem 1). Moreover, in contrast to Yang and Wang (2019), our regret guarantee is anytime, i.e., we don't need to know the value of  $N$  before the algorithm runs.*

**Remark 3** *Considering linear function approximation ( $\mathcal{H}_\varphi = \mathbb{R}^d$ ), Jin et al. (2019) assumes that for any  $V \in \mathcal{V}$  (Equation 3), the map  $(s, a) \mapsto \mathbb{E}_{X \sim P(\cdot|s,a)}[V(X)]$  lies in  $\mathcal{H}_\varphi$ , and propose a model-free algorithm with regret  $\tilde{O}(\sqrt{H^3 d^3 N})$ . For linear kernels, it can be verified that  $\gamma_N = O(d \log N)$  and thus our regret (Theorem 2) is of the order  $\tilde{O}(Hd\sqrt{N})$ . We note that this apparent improvement in our bound is a consequence of slightly stronger assumptions 1 and 2. While they obtain the bound by proving a uniform concentration result over the set  $\mathcal{V}$ , our result uses a novel concentration property of CME estimates (Theorem 1).*

**Proof sketch of Theorem 2** A control on the  $Q$ -function estimates  $Q_h^t$  leads to the regret bound, as our policy is based on  $Q_h^t$ . We prove that as long as  $\Theta_P$  lies in the confidence set  $\mathcal{C}_t$ , the  $Q$ -updates are optimistic estimates of the optimal  $Q$ -values, i.e.,  $Q_h^*(s, a) \leq Q_h^t(s, a)$  for all  $(s, a)$ , and thus, allow us to pick an optimistic action while sufficiently exploring the state space. This implies  $V_1^*(s_1^t) \leq V_1^t(s_1^t)$ , so that the regret  $\mathcal{R}(N) \leq \sum_{t=1}^T (V_1^t(s_1^t) - V_1^{\pi^t}(s_1^t))$ . Letting  $g_1^t(s_1^t) := V_1^t(s_1^t) - V_1^{\pi^t}(s_1^t)$  denote the gap between the most optimistic value and the actual value obtained at episode  $t$ , we then have

$$g_1^t(s_1^t) \leq \sum_{h=1}^H (Q_h^t(s_h^t, a_h^t) - (R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot|s_h^t, a_h^t)} [V_{h+1}^t(X)]) + m_h^t),$$

where  $(m_h^t)_{t,h}$  denotes a martingale difference sequence. We control this via the Azuma-Hoeffding inequality as  $\sum_{t,h} m_{t,h} = O(H\sqrt{N})$ . The rest of the terms inside the summation can be controlled, by Theorem 1 and by design of the confidence set  $\mathcal{C}_t$ , using the predictive variances  $\sigma_{\varphi,t}^2(s_h^t, a_h^t)$ . In fact, for  $\Theta_P \in \mathcal{C}_t$ , it holds that

$$Q_h^t(s, a) - (R(s, a) + \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)]) \leq 2B_V \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi,t}(s, a).$$

Now, the proof can be completed by showing that  $\sum_{t,h} \sigma_{\varphi,t}(s_h^t, a_h^t) = O(\sqrt{HN\gamma_N})$ . Complete proof of this result is given in the supplementary material.

### 4.3. Robustness to model misspecification

Theorem 2 hinges on the fact that any optimistic estimate of the value function can be specified as an element in  $\mathcal{H}_\psi$ . In this section, we study the case when there is a misspecification error. Formally, we consider the following assumption.

**Assumption 3 (Approximate optimistic closure)** *There exists constants  $\zeta > 0$  and  $B_V > 0$ , such that for any  $V \in \mathcal{V}$  (Equation 3), there exists a function  $\tilde{V} \in \mathcal{H}_\psi$  which satisfies  $\|V - \tilde{V}\|_\infty \leq \zeta$  and  $\|\tilde{V}\|_{\mathcal{H}_\psi} \leq B_V$ . We call  $\zeta$  the misspecification error.*

The quality of this approximation will further depend upon how well any  $V \in \mathcal{V}$  can be approximated by a low-norm function in  $\mathcal{H}_\psi$ . One specialization is to the case when  $\mathcal{V} \in \mathcal{C}_b(\mathcal{S})$ , the vector space of continuous and bounded functions on  $\mathcal{S}$ , and  $k_\psi$  is a  $\mathcal{C}_b(\mathcal{S})$ -universal kernel (Steinwart and Christmann, 2008). In this case, we can choose  $\tilde{V}$  such that  $\|V - \tilde{V}\|_\infty$  is arbitrarily small. For technical reasons, we also make the following assumption.

**Assumption 4** *The RKHS  $\mathcal{H}_\psi$  contains the constant functions.*<sup>4</sup>

The following theorem states that our algorithm is in fact robust to a small model misspecification. To achieve this, we only need to adopt a different exploration term in Equation 13 to account for the misspecification error  $\zeta$ . To this end, define the  $Q$ -function updates as

$$Q_h^t(s, a) := R(s, a) + k_{\varphi,t}(s, a)^\top (K_{\varphi,t} + \lambda I)^{-1} v_{h+1}^t + (B_V + \zeta \|1\|_{\mathcal{H}_\psi}) \beta_t(\delta/2) \sigma_{\varphi,t}(s, a), \quad (15)$$

where  $\|1\|_{\mathcal{H}_\psi}$  denotes the norm of the all-one function  $s \mapsto 1$  in  $\mathcal{H}_\psi$ .

**Theorem 3 (Cumulative regret under misspecification)** *Under assumptions 2, 3 and 4, with probability at least  $1 - \delta$ , CME-RL achieves the regret bound*

$$\mathcal{R}(N) \leq 2 (B_V + \zeta \|1\|_{\mathcal{H}_\psi}) \alpha_{N,\delta} \sqrt{2(1 + \lambda^{-1} B_\varphi^2 H) N \gamma_N} + 4\zeta N + 2H \sqrt{2N \log(2/\delta)},$$

4. This is a mild assumption. For any RKHS  $\mathcal{H}_\psi$ , the direct sum  $\mathcal{H}_\psi + \mathbb{R}$ , where  $\mathbb{R}$  denotes the RKHS associated with the kernel  $k(s, s') = 1$ , is again a RKHS with kernel  $k_{\text{new}}(s, s') := k_\psi(s, s') + 1$ .

where  $B_\varphi$  and  $\alpha_{N,\delta}$  are as given in [Theorem 2](#).

In comparison with [Theorem 2](#), [Theorem 3](#) asserts that CME-RL will incur at most an additional  $O(\zeta\gamma_N\sqrt{HN} + \zeta N)$  regret when the model is misspecified. This additional term is linear in  $N$  due to the intrinsic bias introduced by the approximation. This linear dependency is standard in the literature, e.g., it is present even in the easier setting of linear function approximation ([Jin et al., 2019](#)). When  $\zeta$  is sufficiently small (as is typical for universal kernels  $k_\psi$ ), our algorithm will still enjoy good theoretical guarantees.

**Conversion to PAC guarantee** Similar to [Theorem 2](#), we can also convert [Theorem 3](#) to a PAC guarantee. Assuming a fixed initial state  $s$ , with at least a constant probability, we can learn an  $\varepsilon$ -optimal policy  $\pi$  that satisfies  $V_1^*(s) - V_1^\pi(s) \leq \varepsilon + \zeta\gamma_N H^{3/2}$  by running CME-RL for  $T = O(d_{\text{eff}}^2 H^2 / \varepsilon^2)$  episodes.

**Remark 4 (Regret under unknown misspecification error)** *When the misspecification error  $\zeta$  is unknown to the agent a priori, one can invoke the dynamic regret balancing scheme of [Cutkosky et al. \(2021\)](#) to get essentially a similar bound as [Theorem 3](#) (albeit with a polylog factor blow-up). In fact, [Cutkosky et al. \(2021\)](#) gives a bound for the linear MDP model of [Jin et al. \(2019\)](#). Similar techniques can be incorporated to derive a regret bound with unknown  $\zeta$  in our setting also.*

**Proof sketch of [Theorem 3](#)** Similar to the proof of [Theorem 2](#), we control the  $Q$ -function estimates  $Q_h^t(s, a)$  (cf. [Equation 15](#)), but with necessary modifications taking the effect of the misspecification error  $\zeta$  into account. Specifically, we show, for  $\Theta_P \in \mathcal{C}_t$ , that 
$$Q_h^t(s, a) - (R(s, a) + \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)]) \leq 2(B_V + \zeta \|1\|_{\mathcal{H}_\psi}) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) + 2\zeta.$$
 With the result above, we can derive an upper bound on the optimal value  $Q_h^*$  as  $Q_h^*(s, a) \leq Q_h^t(s, a) + 2(H - h)\zeta$ , which allows us to pick an optimistic action. The proof then follows similar steps of [Theorem 2](#) via control of predictive variances and Azuma’s inequality. Complete proof is given in the supplementary material.

## 5. Conclusion

In this paper, we have presented a novel model-based RL algorithm with sub-linear regret guarantees under an optimistic RKHS-closure assumption on the value functions, without requiring a “simulator” access. The algorithm essentially performs an optimistic value iteration step, which is derived from a novel concentration inequality for the mean embeddings of the transition distribution. We have also shown robustness of our algorithm to small model misspecifications.

As future work, it remains an open research direction to relax the strong optimistic closure assumption to a milder one, as in [Zanette et al. \(2020b\)](#) and [Domingues et al. \(2021\)](#), without sacrificing on the computational and regret performances. In terms of computational complexity, [Vial et al. \(2022\)](#) proposed an algorithm for misspecified linear MDPs with bounded per-iteration computational complexity. Although our method has computational complexity growing with the number of data points, we highlight that constant cost per iteration can be achieved with kernel-based approximations by means of low-rank decompositions ([Gijbarts and Metta, 2013](#)), which is another avenue for future work.

## Acknowledgments

The authors would like to thank Sho Takemori and an anonymous reviewer for pointing out a bug in the proof of [Theorem 1](#) in an earlier version of this work.

## References

- Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26, 2011.
- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin F Yang. Model-based reinforcement learning with value-targeted regression. *arXiv preprint arXiv:2006.01107*, 2020.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 844–853. JMLR. org, 2017.
- Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3197–3205, 2019.
- Sayak Ray Chowdhury, Rafael Oliveira, and Fabio Ramos. Active learning of conditional mean embeddings via bayesian optimisation. In *Conference on Uncertainty in Artificial Intelligence*, pages 1119–1128. PMLR, 2020.
- Ashok Cutkosky, Christoph Dann, Abhimanyu Das, Claudio Gentile, Aldo Pacchiano, and Manish Purohit. Dynamic balancing for model selection in bandits and rl. In *International Conference on Machine Learning*, pages 2276–2285. PMLR, 2021.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems 30 (NIPS)*, pages 5711–5721, 2017.
- Aniket Anand Deshmukh, Urun Dogan, and Clay Scott. Multi-task learning for contextual bandits. *Advances in neural information processing systems*, 30, 2017.
- Omar Darwiche Domingues, Pierre Ménard, Matteo Pirotta, Emilie Kaufmann, and Michal Valko. Regret bounds for kernel-based reinforcement learning. *arXiv preprint arXiv:2004.05599*, 2020.
- Omar Darwiche Domingues, Pierre Ménard, Matteo Pirotta, Emilie Kaufmann, and Michal Valko. Kernel-based reinforcement learning: A finite-time analysis. In *International Conference on Machine Learning*, pages 2783–2792. PMLR, 2021.

- Petros Drineas and Michael W Mahoney. On the nyström method for approximating a gram matrix for improved kernel-based learning. *Journal of Machine Learning Research*, 6:2153–2175, 2005.
- Yonathan Efroni, Nadav Merlis, Mohammad Ghavamzadeh, and Shie Mannor. Tight regret bounds for model-based reinforcement learning with greedy policies. In *Advances in Neural Information Processing Systems*, pages 12203–12213, 2019.
- Kenji Fukumizu, Arthur Gretton, Xiaohai Sun, and Bernhard Schölkopf. Kernel measures of conditional dependence. In *Advances in neural information processing systems*, pages 489–496, 2008.
- Kenji Fukumizu, Francis R Bach, Michael I Jordan, et al. Kernel dimension reduction in regression. *The Annals of Statistics*, 37(4):1871–1905, 2009.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 263–272, 2017.
- Arjan Gijsberts and Giorgio Metta. Real-time model learning using Incremental Sparse Spectrum Gaussian Process Regression. *Neural Networks*, 41:59–69, 2013.
- Steffen Grünewälder, Guy Lever, Luca Baldassarre, Massimiliano Pontil, and Arthur Gretton. Modelling transition dynamics in mdps with rkhs embeddings. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 1603–1610, 2012.
- Kelvin Hsu and Fabio Ramos. Bayesian Learning of Conditional Kernel Mean Embeddings for Automatic Likelihood-Free Inference. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, Naha, Okinawa, Japan, 2019.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019.
- Sham Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. *arXiv preprint arXiv:2006.12466*, 2020.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.



- M. N. Lukic and J. H. Beder. Stochastic processes with sample paths in reproducing kernel Hilbert spaces. *Transactions of the American Mathematical Society*, 353(10):3945–3969, 2001.
- Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, and Bernhard Schölkopf. Kernel Mean Embedding of Distributions: A Review and Beyond. *arXiv*, 2016.
- Ronald Ortner and Daniil Ryabko. Online regret bounds for undiscounted continuous reinforcement learning. In *Advances in Neural Information Processing Systems 25 (NIPS)*, pages 1763–1771, 2012.
- Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In *Advances in Neural Information Processing Systems 27 (NIPS)*, pages 1466–1474, 2014.
- Ian Osband, Dan Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems 26 (NIPS)*, pages 3003–3011, 2013.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis, editors, *Proceedings of the 21st conference on advances in Neural Information Processing Systems*, NIPS '07, Vancouver, British Columbia, Canada, dec 2007. MIT Press.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. MIT Press, 2006.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.
- Jonathan Scarlett, Ilija Bogunovic, and Volkan Cevher. Lower bounds on regret for noisy gaussian process bandit optimization. *arXiv preprint arXiv:1706.00090*, 2017.
- Bernhard Schölkopf and Alexander J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT Press, Cambridge, Mass, 2002.
- Le Song, Jonathan Huang, Alex Smola, and Kenji Fukumizu. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 961–968, 2009.
- Le Song, Byron Boots, Sajid M Siddiqi, Geoffrey Gordon, and Alex Smola. Hilbert space embeddings of hidden markov models. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 991–998, 2010a.
- Le Song, Arthur Gretton, and Carlos Guestrin. Nonparametric tree graphical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 765–772, 2010b.

- Le Song, Kenji Fukumizu, and Arthur Gretton. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. *IEEE Signal Processing Magazine*, 30(4):98–111, 2013.
- Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- Bharath K. Sriperumbudur, Kenji Fukumizu, and Gert R. G. Lanckriet. Universality, Characteristic Kernels and RKHS Embedding of Measures. *Journal of Machine Learning Research (JMLR)*, 12:2389–2410, 2011.
- Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- Alexander L Strehl, Lihong Li, and Michael L Littman. Reinforcement learning in finite mdps: PAC analysis. *Journal of Machine Learning Research*, 10(Nov):2413–2444, 2009.
- Daniel Vial, Advait Parulekar, Sanjay Shakkottai, and R Srikant. Improved algorithms for misspecified linear markov decision processes. In *International Conference on Artificial Intelligence and Statistics*, pages 4723–4746. PMLR, 2022.
- Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Provably efficient reinforcement learning with general value function approximation. *arXiv preprint arXiv:2005.10804*, 2020.
- Yining Wang, Ruosong Wang, Simon S Du, and Akshay Krishnamurthy. Optimism in reinforcement learning with generalized linear function approximation. *arXiv preprint arXiv:1912.04136*, 2019.
- Lin F Yang and Mengdi Wang. Reinforcement leaning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019.
- Zhuoran Yang, Chi Jin, Zhaoran Wang, Mengdi Wang, and Michael Jordan. Provably efficient reinforcement learning with kernel and neural function approximations. *Advances in Neural Information Processing Systems*, 33, 2020.
- Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312, 2019.
- Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirota, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964, 2020a.
- Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. *arXiv preprint arXiv:2003.00153*, 2020b.

## Appendix A. Closed-form solution for the operator optimization problem

In this section, we show that Equation 13 is a closed-form solution for the operator optimization problem in Equation 12, which is of the form  $\max_{\Theta \in \mathcal{C}_t} \mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)]$ . Although the term corresponding to  $f$  in Equation 12 is not necessarily in the RKHS  $\mathcal{H}_\psi$ , we may for now assume that  $f \in \mathcal{H}_\psi$ . In this case, the problem above may be rewritten as a quadratically constrained linear program in  $\mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\psi)$ :

$$\max_{\Theta \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\psi)} \langle f, \Theta \varphi(s, a) \rangle_{\mathcal{H}_\psi} \quad \text{s.t.} \quad \left\| (\Theta - \widehat{\Theta}_t) M_t^{1/2} \right\| \leq \beta_t(\delta/2), \quad (16)$$

where  $M_t \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\varphi)$ . The problem above admits a closed-form solution by applying the Karush-Kuhn-Tucker (KKT) conditions. Directly doing so would require us to take derivatives of the operator norm. However, observe that for any operator  $M \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\varphi)$  we have  $\|M\| \leq \sqrt{\text{tr}(M^\top M)}$ , where the latter corresponds to the Hilbert-Schmidt norm. Compared to the operator norm, we can easily take derivatives of the trace to compute the KKT conditions. In addition, due to the upper bound, any solution satisfying the Hilbert-Schmidt norm constraint is also a solution under the operator norm constraint. We replace Equation 16 with the following problem:

$$\max_{\Theta \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\psi)} \langle f, \Theta \varphi(s, a) \rangle_{\mathcal{H}_\psi} \quad \text{s.t.} \quad \text{tr}((\Theta - \widehat{\Theta}_t) M_t (\Theta - \widehat{\Theta}_t)^\top) \leq \beta_t(\delta/2)^2.$$

Applying the KKT conditions, we solve  $\nabla_{\Theta} \ell(\Theta, \eta) = 0$  with respect to  $\Theta \in \mathcal{L}(\mathcal{H}_\varphi, \mathcal{H}_\psi)$  and  $\eta \in \mathbb{R}$ ,  $\eta \geq 0$ , where  $\ell(\Theta, \eta) := \langle f, \Theta \varphi(s, a) \rangle_{\mathcal{H}_\psi} - \eta(\text{tr}((\Theta - \widehat{\Theta}_t) M_t (\Theta - \widehat{\Theta}_t)^\top) - \beta_t(\delta/2)^2)$ . First, by  $\nabla_{\Theta} \ell(\Theta, \eta) = 0$ , we have

$$f \otimes \varphi(s, a) - 2\eta(\Theta - \widehat{\Theta}_t) M_t = 0 \implies \Theta = \widehat{\Theta}_t + \frac{1}{2\eta} (f \otimes \varphi(s, a)) M_t^{-1}. \quad (17)$$

Now, note that, for quadratically constrained linear program, the maximum should lie at the border of the constrained set. Replacing the result above into the constraint, we obtain

$$\beta_t(\delta/2)^2 = \frac{1}{4\eta^2} \text{tr}((f \otimes \varphi(s, a)) M_t^{-1} (\varphi(s, a) \otimes f)) = \frac{1}{4\eta^2} \langle \varphi(s, a), M_t^{-1} \varphi(s, a) \rangle_{\mathcal{H}_\varphi} \langle f, f \rangle_{\mathcal{H}_\psi},$$

so that  $\eta = \frac{1}{2\beta_t(\delta/2)} \|\varphi(s, a)\|_{M_t^{-1}} \|f\|_{\mathcal{H}_\psi}$ . Combining the latter with Equation 17, the solution to Equation 16 is then given by

$$\Theta_* := \widehat{\Theta}_t + \frac{\beta_t(\delta/2)}{\|\varphi(s, a)\|_{M_t^{-1}} \|f\|_{\mathcal{H}_\psi}} (f \otimes \varphi(s, a)) M_t^{-1},$$

which finally yields

$$\begin{aligned} \max_{\Theta \in \mathcal{C}_t} \mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)] &= \langle f, \Theta_* \varphi(s, a) \rangle_{\mathcal{H}_\psi} = \left\langle f, \widehat{\Theta}_t \varphi(s, a) \right\rangle_{\mathcal{H}_\psi} + \beta_t(\delta/2) \|f\|_{\mathcal{H}_\psi} \|\varphi(s, a)\|_{M_t^{-1}} \\ &= \left\langle f, \widehat{\Theta}_t \varphi(s, a) \right\rangle_{\mathcal{H}_\psi} + \beta_t(\delta/2) \lambda^{-1/2} \|f\|_{\mathcal{H}_\psi} \sigma_{\varphi, t}(s, a). \end{aligned}$$

Replacing  $f$  by  $V_{h+1}^t$  in the solution above and adding the reward function (cf. Equation 12), we recover Equation 13.

## Appendix B. Proof of main results

### B.1. Proof of Theorem 1

We first need to define the data-generating process and its properties. Let  $\mathcal{F}_{t,h-1}$  be the filtration induced by the sequence  $\mathcal{D}_t \cup \{(s_j^t, a_j^t)\}_{j \leq h}$ , where  $\mathcal{D}_t$  denotes the replay buffer at the beginning of episode  $t$ . Note that  $\mathcal{F}_{t,0} = \mathcal{D}_t$ . The evaluation noise defined by  $\varepsilon_h^t := \psi(s_{h+1}^t) - \Theta_P \varphi(s_h^t, a_h^t)$ , for  $s_{h+1}^t \sim P(\cdot | s_h^t, a_h^t)$ , is such that:

$$\mathbb{E}[\varepsilon_h^t | \mathcal{F}_{t,h-1}] = 0 \quad \text{and} \quad \|\varepsilon_h^t\|_{\mathcal{H}_\psi} \leq 2 \quad \text{a.s.}$$

We overload notation to define, for each pair  $(t, h)$ , the operator  $M_{t,h} = M_t + \sum_{j \leq h} \varphi(s_j^t, a_j^t) \otimes \varphi(s_j^t, a_j^t)$ , where  $M_t = \sum_{\tau < t, h \leq H} \varphi(s_h^\tau, a_h^\tau) \otimes \varphi(s_h^\tau, a_h^\tau) + \lambda I$ . Note that  $M_{t,h} = M_{t,h-1} + \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t)$ , where  $M_{t,0} = M_t = M_{t-1,H}$ . Define the estimate

$$\widehat{\Theta}_{t,h} = \left( \sum_{\tau < t, h \leq H} \psi(s_{h+1}^\tau) \otimes \varphi(s_h^\tau, a_h^\tau) + \sum_{j \leq h} \psi(s_{j+1}^t) \otimes \varphi(s_j^t, a_j^t) \right) M_{t,h}^{-1}.$$

Note that  $\widehat{\Theta}_t = \widehat{\Theta}_{t-1,H}$ , where  $\widehat{\Theta}_t$  is given by (7). Consider the random variable

$$z_{t,h} = \left\| (\widehat{\Theta}_{t,h} - \Theta_P) M_{t,h}^{1/2} \right\|_{\text{HS}}^2 = \text{tr} \left( (\widehat{\Theta}_{t,h} - \Theta_P) M_{t,h} (\widehat{\Theta}_{t,h} - \Theta_P)^\top \right).$$

Define the operator  $X_{t,h} = (\widehat{\Theta}_{t,h} - \Theta_P) M_{t,h}$ . Note that

$$X_{t,h} = \sum_{\tau < t, h \leq H} \varepsilon_h^\tau \otimes \varphi(s_h^\tau, a_h^\tau) + \sum_{j \leq h} \varepsilon_j^t \otimes \varphi(s_j^t, a_j^t) - \lambda \Theta_P = X_{t,h-1} + \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t),$$

where  $X_{t,0} = X_{t-1,H}$ . We then have

$$\begin{aligned} z_{t,h} &= \text{tr} \left( X_{t,h} M_{t,h}^{-1} X_{t,h}^\top \right) \\ &= \underbrace{\text{tr} \left( X_{t,h-1} M_{t,h}^{-1} X_{t,h-1}^\top \right)}_A + 2 \underbrace{\text{tr} \left( \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) M_{t,h}^{-1} X_{t,h-1}^\top \right)}_B + \underbrace{\text{tr} \left( \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) M_{t,h}^{-1} \varphi(s_h^t, a_h^t) \otimes \varepsilon_h^t \right)}_C. \end{aligned}$$

Define the variance estimate  $\sigma_{\varphi,t,h}^2(s, a) = \lambda \left\langle \varphi(s, a), M_{t,h}^{-1} \varphi(s, a) \right\rangle_{\mathcal{H}_\varphi}$ . Note that  $\sigma_{\varphi,t}^2(s, a) = \sigma_{\varphi,t-1,H}^2(s, a)$ . Now, the Sherman-Morrison formula yields

$$\begin{aligned} M_{t,h}^{-1} &= (M_{t,h-1} + \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t))^{-1} = M_{t,h-1}^{-1} - \frac{M_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t) M_{t,h-1}^{-1}}{1 + \left\langle \varphi(s_h^t, a_h^t), M_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \right\rangle_{\mathcal{H}_\varphi}} \\ &= M_{t,h-1}^{-1} - \frac{M_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t) M_{t,h-1}^{-1}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}. \end{aligned}$$

We thus have

$$\begin{aligned}
 \text{A} &= \text{tr} \left( \mathbf{X}_{t,h-1} \mathbf{M}_{t,h-1}^{-1} \mathbf{X}_{t,h-1}^\top \right) - \text{tr} \left( \mathbf{X}_{t,h-1} \frac{\mathbf{M}_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t) \mathbf{M}_{t,h-1}^{-1}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \mathbf{X}_{t,h-1}^\top \right) \\
 &\leq \text{tr} \left( \mathbf{X}_{t,h-1} \mathbf{M}_{t,h-1}^{-1} \mathbf{X}_{t,h-1}^\top \right) = z_{t,h-1}, \\
 \text{B} &= \text{tr} \left( \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) \mathbf{M}_{t,h-1}^{-1} \mathbf{X}_{t,h-1}^\top \right. \\
 &\quad \left. - \frac{\varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) \mathbf{M}_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t) \mathbf{M}_{t,h-1}^{-1} \mathbf{X}_{t,h-1}^\top}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \right) \\
 &= \text{tr} \left( \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) (\widehat{\Theta}_{t,h-1} - \Theta_P)^\top \right. \\
 &\quad \left. - \frac{\left\langle \varphi(s_h^t, a_h^t), \mathbf{M}_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \right\rangle_{\mathcal{H}_\varphi} \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) (\widehat{\Theta}_{t,h-1} - \Theta_P)^\top}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \right) \\
 &= \left( 1 - \frac{\lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \right) \text{tr} \left( \varepsilon_h^t \otimes \varphi(s_h^t, a_h^t) (\widehat{\Theta}_{t,h-1} - \Theta_P)^\top \right) \\
 &= \frac{\left\langle (\widehat{\Theta}_{t,h-1} - \Theta_P) \varphi(s_h^t, a_h^t), \varepsilon_h^t \right\rangle_{\mathcal{H}_\psi}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}, \quad \text{and} \\
 \text{C} &= \|\varepsilon_h^t\|_{\mathcal{H}_\psi}^2 \left( \left\langle \varphi(s_h^t, a_h^t), \mathbf{M}_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \right\rangle_{\mathcal{H}_\varphi} \right. \\
 &\quad \left. - \left\langle \varphi(s_h^t, a_h^t), \frac{\mathbf{M}_{t,h-1}^{-1} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t) \mathbf{M}_{t,h-1}^{-1}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \varphi(s_h^t, a_h^t) \right\rangle_{\mathcal{H}_\varphi} \right) \\
 &= \|\varepsilon_h^t\|_{\mathcal{H}_\psi}^2 \left( \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t) - \frac{\lambda^{-2} \sigma_{\varphi,t,h-1}^4(s_h^t, a_h^t)}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \right) \\
 &= \|\varepsilon_h^t\|_{\mathcal{H}_\psi}^2 \frac{\lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}.
 \end{aligned}$$

Putting these together, we have

$$\begin{aligned}
 z_{t,h} &\leq z_{t,h-1} + 2 \frac{\left\langle (\widehat{\Theta}_{t,h-1} - \Theta_P) \varphi(s_h^t, a_h^t), \varepsilon_h^t \right\rangle_{\mathcal{H}_\psi}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} + \|\varepsilon_h^t\|_{\mathcal{H}_\psi}^2 \frac{\lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \\
 &\leq \lambda \|\Theta_P\|_{\text{HS}}^2 + 2 \sum_{(\tau,j) \leq (t,h)} \frac{\left\langle (\widehat{\Theta}_{\tau,j-1} - \Theta_P) \varphi(s_j^\tau, a_j^\tau), \varepsilon_j^\tau \right\rangle_{\mathcal{H}_\psi}}{1 + \lambda^{-1} \sigma_{\varphi,\tau,j-1}^2(s_j^\tau, a_j^\tau)} \\
 &\quad + \sum_{(\tau,j) \leq (t,h)} \|\varepsilon_j^\tau\|_{\mathcal{H}_\psi}^2 \frac{\lambda^{-1} \sigma_{\varphi,\tau,j-1}^2(s_j^\tau, a_j^\tau)}{1 + \lambda^{-1} \sigma_{\varphi,\tau,j-1}^2(s_j^\tau, a_j^\tau)}.
 \end{aligned}$$

For any  $B_P \geq \|\Theta_P\|_{\text{HS}}$ , define

$$\beta_{t,h}(\delta) := \sqrt{2\lambda B_P^2 + 256(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,t,h})^{1/2}) \log(2t^2 H/\delta)},$$

where  $\mathbf{K}_{\varphi,t,h} := [k_\varphi((s_j^\tau, a_j^\tau), (s_{j'}^{\tau'}, a_{j'}^{\tau'}))]_{(\tau,j),(\tau',j') \leq (t,h)}$  denotes the gram-matrix at step  $h$  of episode  $t$ . Note that  $\mathbf{K}_{\varphi,t} = \mathbf{K}_{\varphi,t-1,H}$ . Now define an event  $\mathcal{E}_{t,h}$  as

$$\mathcal{E}_{t,h} = \mathbb{I}\{z_{\tau,j} \leq \beta_{\tau,j}^2(\delta), \quad \forall (\tau,j) \leq (t,h)\}.$$

Note that  $\mathcal{E}_{t,h} = 1$  implies  $\mathcal{E}_{t,h-1} = 1$  for all  $(t,h) \geq (1,1)$ , where  $\mathcal{E}_{t,0} = \mathcal{E}_{t-1,H}$ . Now we

define a sequence of random variables  $\{y_{t,h}\}_{t,h}$  as  $y_{t,h} = \mathcal{E}_{t,h-1} \frac{\langle (\widehat{\Theta}_{t,h-1} - \Theta_P) \varphi(s_h^t, a_h^t), \varepsilon_h^t \rangle_{\mathcal{H}_\psi}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)}$ .

Note that  $y_{t,h}$  is  $\mathcal{F}_{t,h}$ -measurable and  $\mathbb{E}[y_{t,h} | \mathcal{F}_{t,h-1}] = 0$ . Hence  $\{y_{t,h}\}_{t,h}$  is a martingale difference sequence w.r.t. the filtration  $\{\mathcal{F}_{t,h}\}_{t,h}$ . Note that

$$\begin{aligned}
 |y_{t,h}| &\leq \mathcal{E}_{t,h-1} \|\varepsilon_h^t\|_{\mathcal{H}_\psi} \frac{\left\| (\widehat{\Theta}_{t,h-1} - \Theta_P) \mathbf{M}_{t,h-1}^{1/2} \right\|_{\text{HS}} \left\| \mathbf{M}_{t,h-1}^{-1/2} \varphi(s_h^t, a_h^t) \right\|_{\mathcal{H}_\varphi}}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \\
 &\leq \beta_{t,h-1}(\delta) \|\varepsilon_h^t\|_{\mathcal{H}_\psi} \frac{\lambda^{-1} \sigma_{\varphi,t,h-1}(s_h^t, a_h^t)}{1 + \lambda^{-1} \sigma_{\varphi,t,h-1}^2(s_h^t, a_h^t)} \leq 2\beta_{t,h-1}(\delta) \lambda^{-1} \sigma_{\varphi,t,h-1}(s_h^t, a_h^t),
 \end{aligned}$$

since  $\|\varepsilon_h^t\|_{\mathcal{H}_\psi} \leq 2$  a.s. We then have

$$\begin{aligned}
 \sum_{(\tau,j) \leq (t,h)} |y_{\tau,j}|^2 &\leq \sum_{(\tau,j) \leq (t,h)} 4\beta_{\tau,j-1}^2(\delta) \lambda^{-1} \sigma_{\varphi,\tau,j-1}^2(s_j^\tau, a_j^\tau) \\
 &\leq 8(1 + \lambda^{-1}) \beta_{t,h}^2(\delta) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,t,h})^{1/2}),
 \end{aligned}$$

where we have used that  $\sum_{(\tau,j) \leq (t,h)} \lambda^{-1} \sigma_{\varphi,\tau,j-1}^2(s_j^\tau, a_j^\tau) \leq 2(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,t,h})^{1/2})$ .

Then, by Azuma-Hoeffding's inequality, with probability at least  $1 - \frac{\delta}{2t^2 H}$ , we have

$$\sum_{(\tau,j) \leq (t,h)} y_{\tau,j} \leq \sqrt{16(1 + \lambda^{-1}) \beta_{t,h}^2(\delta) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,t,h})^{1/2}) \log(2t^2 H/\delta)} \leq \beta_{t,h}^2(\delta)/4,$$

as  $\beta_{t,h}^2(\delta) \geq 256(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi,t,h})^{1/2}) \log(2t^2 H/\delta)$ . Hence, by an union bound

$$\mathbb{P}\left[\exists (t,h) \geq (1,1) : \sum_{(\tau,j) \leq (t,h)} y_{\tau,j} > \beta_{t,h}^2(\delta)/4\right] \leq \sum_{t=1}^{\infty} \sum_{h=1}^H \frac{\delta}{2t^2 H} \leq \frac{\delta \pi^2}{12} \leq \delta. \quad (18)$$

Now, it suffices to show that  $z_{t,h} \leq \beta_{t,h}^2(\delta)$  for all  $(t,h) \geq (1,0)$  given  $\sum_{(\tau,j) \leq (t,h)} y_{\tau,j} \leq \beta_{t,h}^2(\delta)/4$ , for all  $(t,h) \geq (1,1)$ . We will show this by induction on  $(t,h)$ . For the base case  $(t,h) = (1,0)$ , we have  $z_{1,0} = \lambda \|\Theta_P\|_{\text{HS}}^2 \leq \beta_{1,0}^2(\delta)$ . Now by inductive hypothesis, let



$z_{\tau,j} \leq \beta_{\tau,j}^2(\delta)$  for all  $(1,0) \leq (\tau,j) \leq (t,h-1)$ . We then have  $\mathcal{E}_{\tau,j} = 1$  for all  $(1,0) \leq (\tau,j) \leq (t,h-1)$ . Therefore, we have

$$\begin{aligned} z_{t,h} &\leq \lambda \|\Theta_P\|_{\text{HS}}^2 + 2 \sum_{(\tau,j) \leq (t,h)} \mathcal{E}_{\tau,j-1} \frac{\langle (\widehat{\Theta}_{\tau,j-1} - \Theta_P) \varphi(s_j^\tau, a_j^\tau), \varepsilon_j^\tau \rangle_{\mathcal{H}_\psi}}{1 + \lambda^{-1} \sigma_{\varphi, \tau, j-1}^2(s_j^\tau, a_j^\tau)} \\ &\quad + \sum_{(\tau,j) \leq (t,h)} \|\varepsilon_j^\tau\|_{\mathcal{H}_\psi}^2 \frac{\lambda^{-1} \sigma_{\varphi, \tau, j-1}^2(s_j^\tau, a_j^\tau)}{1 + \lambda^{-1} \sigma_{\varphi, \tau, j-1}^2(s_j^\tau, a_j^\tau)} \\ &\leq \lambda \|\Theta_P\|_{\text{HS}}^2 + 2 \sum_{(\tau,j) \leq (t,h)} y_{\tau,j} + 4 \sum_{(\tau,j) \leq (t,h)} \frac{\lambda^{-1} \sigma_{\varphi, \tau, j-1}^2(s_j^\tau, a_j^\tau)}{1 + \lambda^{-1} \sigma_{\varphi, \tau, j-1}^2(s_j^\tau, a_j^\tau)} \\ &\leq \lambda \|\Theta_P\|_{\text{HS}}^2 + \beta_{t,h}^2(\delta)/2 + 8(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi, t, h})^{1/2}) \leq \beta_{t,h}^2(\delta), \end{aligned}$$

as  $\beta_{t,h}^2(\delta) \geq 2\lambda \|\Theta_P\|_{\text{HS}}^2 + 16(1 + \lambda^{-1}) \log(\det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi, t, h})^{1/2})$ . Now the proof follows from Equation 18 and noting that  $\|(\widehat{\Theta}_t - \Theta_P) \mathbf{M}_t^{1/2}\|_{\text{HS}} = z_{t-1, H}^{1/2} \leq \beta_{t-1, H}(\delta) \leq \beta_t(\delta)$ .

## B.2. Regret analysis of CME-RL

To prove the regret bound in Theorem 2, we establish a sequence of intermediate results to bound the performance gap between the optimal policy and the policy followed by CME-RL. In Lemma 1, we start with a control on the  $Q$ -function estimates  $Q_h^t$ , which in turn leads to the regret bound, as our policy is based on  $Q_h^t$ . The result implies that as long as the true transition distribution lies in the confidence set  $\mathcal{C}_t$ , the  $Q$ -updates are optimistic estimates of the optimal  $Q$ -values and thus, allow us to pick an optimistic action while sufficiently exploring the state space.

**Lemma 1 (Optimism)** *Let  $P \in \mathcal{C}_t$ . Then,  $Q_h^*(s, a) \leq Q_h^t(s, a)$  for all  $h$ , and  $(s, a)$ .*

**Proof** We prove the lemma by induction on  $h$ . When  $h = H$ , the inequality holds by definition. Now we assume that the lemma holds for some  $h' = h + 1$ , where  $1 \leq h < H$ . This implies that for all  $s \in \mathcal{S}$ ,

$$V_{h+1}^t(s) = \min \left\{ H, \max_{a \in \mathcal{A}} Q_{h+1}^t(s, a) \right\} \geq \min \left\{ H, \max_{a \in \mathcal{A}} Q_{h+1}^*(s, a) \right\} = V_{h+1}^*(s).$$

We then have, for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , that

$$\begin{aligned} Q_h^*(s, a) &= R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^*(X)] \\ &\leq R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^t(X)] \\ &\leq R(s, a) + \max_{\Theta_{P'} \in \mathcal{C}_t} \mathbb{E}_{X \sim P'(\cdot | s, a)} [V_{h+1}^*(X)] \leq Q_h^t(s, a), \end{aligned}$$

where the third step follows from  $\Theta_P \in \mathcal{C}_t$ . ■

**Lemma 2 (Gap between optimistic and actual values)** *Let  $g_h^t(s) = V_h^t(s) - V_h^{\pi^t}(s)$ , and  $m_h^t = \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [g_{h+1}^t(X)] - g_{h+1}^t(s_{h+1}^t)$ . Then*

$$g_1^t(s_1^t) \leq \sum_{h=1}^H Q_h^t(s_h^t, a_h^t) - \left( R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)] \right) + \sum_{h=1}^H m_h^t.$$

**Proof** Note that  $a_h^t = \pi_h^t(s_h^t) = \operatorname{argmax}_{a \in \mathcal{A}} Q_h^t(s_h^t, a)$ . Therefore

$$\begin{aligned} V_h^{\pi^t}(s_h^t) &= R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^{\pi^t}(X)] \\ &= R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)] - \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [g_{h+1}^t(X)] \\ &= R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)] - g_{h+1}^t(s_{h+1}^t) - m_h^t. \end{aligned}$$

We also have  $V_h^t(s_h^t) = \min \{H, \max_{a \in \mathcal{A}} Q_h^t(s_h^t, a)\} = \min \{H, Q_h^t(s_h^t, a_h^t)\} \leq Q_h^t(s_h^t, a_h^t)$ . Therefore,  $g_h^t(s_h^t) \leq Q_h^t(s_h^t, a_h^t) - \left(R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)]\right) + g_{h+1}^t(s_{h+1}^t) + m_h^t$ . Since  $g_{H+1}^t(s) = 0$  for all  $s \in \mathcal{S}$ , a simple recursion over all  $h \in [H]$  completes the proof.  $\blacksquare$

**Lemma 3 (Cumulative regret expressed through  $Q$ -estimates)** *Let  $\Theta_P \in \mathcal{C}_t$  for all  $t \geq 1$ . Then for any  $\delta \in (0, 1]$ , the following holds with probability at least  $1 - \delta/2$ :*

$$\mathcal{R}(N) \leq \sum_{t \leq T, h \leq H} Q_h^t(s_h^t, a_h^t) - \left(R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)]\right) + 2H \sqrt{2N \log(2/\delta)}.$$

**Proof** If  $\Theta_P \in \mathcal{C}_t$  for all  $t$ , we have from Lemma 1 that

$$\forall t \geq 1, \quad V_1^t(s_1^t) = \min \left\{ H, \max_{a \in \mathcal{A}} Q_1^t(s_1^t, a) \right\} \geq \min \left\{ H, \max_{a \in \mathcal{A}} Q_1^*(s_1^t, a) \right\} = V_1^*(s_1^t).$$

Therefore the cumulative regret after  $N = TH$  steps is given by

$$\mathcal{R}(N) = \sum_{t=1}^T [V_1^*(s_1^t) - V_1^{\pi^t}(s_1^t)] \leq \sum_{t=1}^T (V_1^t(s_1^t) - V_1^{\pi^t}(s_1^t)) = \sum_{t=1}^T g_1^t(s_1^t).$$

We then have from Lemma 2 that

$$\mathcal{R}(N) \leq \sum_{t \leq T, h \leq H} Q_h^t(s_h^t, a_h^t) - \left(R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)]\right) + \sum_{t \leq T, h \leq H} m_h^t.$$

Note that  $(m_h^t)_{t,h}$  is a martingale difference sequence adapted to the filtration  $\mathcal{F}_{t,h}$  with  $|m_h^t| \leq 2H$ . Hence, by Azuma-Hoeffding inequality, with probability at least  $1 - \delta/2$ ,  $\sum_{t \leq T, h \leq H} m_h^t \leq 2H \sqrt{2TH \log(2/\delta)} = 2H \sqrt{2N \log(2/\delta)}$ , which proves the result.  $\blacksquare$

**Lemma 4 (Error in  $Q$ -estimates)** *Let  $\Theta_P \in \mathcal{C}_t$ . Then, for all  $h \leq H$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,*

$$Q_h^t(s, a) - \left(R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^t(X)]\right) \leq 2B_V \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi, t}(s, a).$$

**Proof** For  $\Theta_P \in \mathcal{C}_t$ , uniformly over all  $t \in \mathbb{N}$ , the mean embeddings satisfy

$$\begin{aligned} \left\| \vartheta_P^{(s,a)} - \widehat{\vartheta}_t^{(s,a)} \right\|_{\mathcal{H}_\psi} &= \left\| (\Theta_P - \widehat{\Theta}_t) \varphi(s, a) \right\|_{\mathcal{H}_\psi} \\ &\leq \left\| (\Theta_P - \widehat{\Theta}_t) (\widehat{\mathbb{C}}_{\varphi, t} + \lambda \mathbf{I})^{1/2} \right\| \left\| (\widehat{\mathbb{C}}_{\varphi, t} + \lambda \mathbf{I})^{-1/2} \varphi(s, a) \right\|_{\mathcal{H}_\varphi} \\ &\leq \beta_t(\delta/2) \left\langle \varphi(s, a), (\widehat{\mathbb{C}}_{\varphi, t} + \lambda \mathbf{I})^{-1} \varphi(s, a) \right\rangle_{\mathcal{H}_\varphi}^{1/2} = \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi, t}(s, a). \end{aligned}$$

Now, by Assumption 1,  $V_{h+1}^t \in \mathcal{H}_\psi$ . Hence, the  $Q$ -estimates (Equation 13) can be written as  $Q_h^t(s, a) = R(s, a) + \left\langle V_{h+1}^t, \widehat{\vartheta}_t^{(s,a)} \right\rangle_{\mathcal{H}_\psi} + B_V \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi, t}(s, a)$ . Therefore, we have

$$\begin{aligned} &Q_h^t(s, a) - \left(R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^t(X)]\right) \\ &= \left\langle V_{h+1}^t, \widehat{\vartheta}_t^{(s,a)} - \vartheta_P^{(s,a)} \right\rangle_{\mathcal{H}_\psi} + B_V \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi, t}(s, a) \leq 2B_V \beta_t(\delta/2) \lambda^{-1/2} \sigma_{\varphi, t}(s, a), \end{aligned}$$

where the last step holds since  $\|V_{h+1}^t\|_{\mathcal{H}_\psi} \leq B_V$  and  $\Theta_P \in \mathcal{C}_t$ .  $\blacksquare$

**Lemma 5 (Sum of predictive variances)** *Let  $\sup_{s,a} \sqrt{k((s,a), (s,a))} \leq B_\varphi$ . Then*

$$\sum_{t \leq T, h \leq H} \lambda^{-1} \sigma_{\varphi,t}^2(s_h^t, a_h^t) \leq (1 + \lambda^{-1} B_\varphi^2 H) \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi, T+1}) .$$

**Proof** We have from Equation 9 that  $\lambda^{-1} \sigma_{\varphi,t}^2(s, a) = \text{tr}(\mathbf{M}_t^{-1} \varphi(s, a) \otimes \varphi(s, a))$ . We also note that  $\mathbf{M}_t \geq \lambda \mathbf{I}$  and  $\varphi(s, a) \otimes \varphi(s, a) \leq B_\varphi^2 \mathbf{I}$ . Therefore  $\mathbf{M}_t^{-1} \varphi(s, a) \otimes \varphi(s, a) \leq \lambda^{-1} B_\varphi^2 \mathbf{I}$ . Now, since  $\mathbf{M}_{t+1} = \mathbf{M}_t + \sum_{h \leq H} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t)$ , we have

$$\mathbf{M}_t^{-1} = (\mathbf{I} + \mathbf{M}_t^{-1} \sum_{h \leq H} \varphi(s_h^t, a_h^t) \otimes \varphi(s_h^t, a_h^t)) \mathbf{M}_{t+1}^{-1} \preceq (1 + \lambda^{-1} B_\varphi^2 H) \mathbf{M}_{t+1}^{-1} .$$

We then have

$$\begin{aligned} \sum_{t \leq T, h \leq H} \lambda^{-1} \sigma_{\varphi,t}^2(s_h^t, a_h^t) &\leq (1 + \lambda^{-1} B_\varphi^2 H) \sum_{t \leq T} \text{tr}(\mathbf{M}_{t+1}^{-1} (\mathbf{M}_{t+1} - \mathbf{M}_t)) \\ &\leq (1 + \lambda^{-1} B_\varphi^2 H) \sum_{t \leq T} \log \frac{\det(\mathbf{M}_{t+1})}{\det(\mathbf{M}_t)} \\ &= (1 + \lambda^{-1} B_\varphi^2 H) \log \det(\mathbf{I} + \lambda^{-1} \mathbf{M}_{T+1}) \\ &= (1 + \lambda^{-1} B_\varphi^2 H) \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{\varphi, T+1}) . \end{aligned}$$

Here, in the second step we have used that for two positive definite operators  $\mathbf{A}$  and  $\mathbf{B}$  such that  $\mathbf{A} - \mathbf{B}$  is positive semi-definite,  $\text{tr}(\mathbf{A}^{-1}(\mathbf{A} - \mathbf{B})) \leq \log \frac{\det(\mathbf{A})}{\det(\mathbf{B})}$ . The last step follows from Sylvester's determinant identity.  $\blacksquare$

Based on the results in this section, we can finally derive a proof for Theorem 2.

### B.2.1. PROOF OF THEOREM 2

We have from Lemma 3 and 4 that if  $\Theta_P \in \mathcal{C}_t$  for all  $t$ , then with probability at least  $1 - \delta/2$ , the cumulative regret

$$\begin{aligned} \mathcal{R}(N) &\leq 2B_V \beta_T(\delta/2) \sum_{t \leq T, h \leq H} \lambda^{-1/2} \sigma_{\varphi,t}(s_h^t, a_h^t) + 2H \sqrt{2N \log(2/\delta)} \\ &\leq 2B_V \alpha_{N,\delta} \sqrt{TH \sum_{t \leq T, h \leq H} \lambda^{-1} \sigma_{\varphi,t}^2(s_h^t, a_h^t)} + 2H \sqrt{2N \log(2/\delta)} \\ &\leq 2B_V \alpha_{N,\delta} \sqrt{2(1 + \lambda^{-1} B_\varphi^2 H) N \gamma_N} + 2H \sqrt{2N \log(2/\delta)} . \end{aligned}$$

The first step follows since  $\beta_t(\delta)$  increases with  $t$ , the second step is due to Cauchy-Schwartz's inequality and the fact that  $\beta_T(\delta/2) \leq \alpha_{N,\delta}$ , and the final step follows from Lemma 5. The proof now can be completed using Theorem 1 and taking a union bound.

### B.3. Regret analysis of CME-RL under model misspecification

To prove the regret bound in Theorem 3, we follow the similar arguments used in proving Theorem 2, but with necessary modifications taking the effect of the misspecification error  $\zeta$  into account. We first derive the following result.

**Lemma 6 (Error in approximate Q-values)** *Let  $\Theta_P \in \mathcal{C}_t$ . Then, for all  $h \in [H]$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we have that:*

$$Q_h^t(s, a) - (R(s, a) + \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)]) \leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) + 2\zeta .$$

**Proof** Note that the Q-estimates (Equation 15) can be rewritten as

$$Q_h^t(s, a) := R(s, a) + \alpha_t(s, a)^\top v_{h+1}^t + \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) .$$

By Assumption 3, there exists a function  $\tilde{V}_{h+1}^t \in \mathcal{H}_\psi$  such that  $\|V_{h+1}^t - \tilde{V}_{h+1}^t\|_\infty \leq \zeta$ . We now define the vector  $\tilde{v}_{h+1}^t := [\tilde{V}_{h+1}^t(s_{h'+1}^\tau)]_{\tau < t, h' \leq H}$  and introduce the shorthand notation  $\mathbb{E}_{P(\cdot|s,a)}[f] := \mathbb{E}_{X \sim P(\cdot|s,a)}[f(X)]$ . We then have

$$\begin{aligned} & \left| \alpha_t(s, a)^\top v_{h+1}^t - \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)] \right| \\ &= \left| \alpha_t(s, a)^\top \tilde{v}_{h+1}^t + \alpha_t(s, a)^\top (v_{h+1}^t - \tilde{v}_{h+1}^t) - \mathbb{E}_{P(\cdot|s,a)}[\tilde{V}_{h+1}^t] + \mathbb{E}_{P(\cdot|s,a)}[\tilde{V}_{h+1}^t - V_{h+1}^t] \right| \\ &\leq \left| \alpha_t(s, a)^\top \tilde{v}_{h+1}^t - \mathbb{E}_{P(\cdot|s,a)}[\tilde{V}_{h+1}^t] \right| + \|\alpha_t(s, a)\|_1 \left\| v_{h+1}^t - \tilde{v}_{h+1}^t \right\|_\infty + \left\| \tilde{V}_{h+1}^t - V_{h+1}^t \right\|_\infty \\ &\leq \left| \alpha_t(s, a)^\top \tilde{v}_{h+1}^t - \mathbb{E}_{P(\cdot|s,a)}[\tilde{V}_{h+1}^t] \right| + \zeta (1 + \|\alpha_t(s, a)\|_1) , \end{aligned} \quad (19)$$

which follows by an application of Hölder's inequality. Now, as  $\Theta_P \in \mathcal{C}_t$  and  $\left\| \tilde{V}_{h+1}^t \right\|_{\mathcal{H}_\psi} \leq B_V$ , the following also holds:

$$\begin{aligned} \left| \alpha_t(s, a)^\top \tilde{v}_{h+1}^t - \mathbb{E}_{P(\cdot|s,a)}[\tilde{V}_{h+1}^t] \right| &= \left| \left\langle \tilde{V}_{h+1}^t, \hat{\vartheta}_t(s, a) - \vartheta_P(s, a) \right\rangle_{\mathcal{H}_\psi} \right| \\ &\leq B_V \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) . \end{aligned} \quad (20)$$

By Assumption 4, the constant function  $1 : \mathcal{S} \rightarrow \mathbb{R}$  is an element of  $\mathcal{H}_\psi$ . For stationary (radial) kernels  $k_\varphi$ , we have  $[\alpha_t(s, a)]_{\tau, h} \geq 0, \forall \tau \leq t, h \leq H$ . Now, as  $\Theta_P \in \mathcal{C}_t$ , we have:

$$\begin{aligned} \|\alpha_t(s, a)\|_1 &= \left\langle \hat{\vartheta}_t^{(s,a)}, 1 \right\rangle_{\mathcal{H}_\psi} = \left\langle \vartheta_P^{(s,a)}, 1 \right\rangle_{\mathcal{H}_\psi} + \left\langle \hat{\vartheta}_t^{(s,a)} - \vartheta_P^{(s,a)}, 1 \right\rangle_{\mathcal{H}_\psi} \\ &\leq \mathbb{E}_{X \sim P(\cdot|s,a)}[1(X)] + \|1\|_{\mathcal{H}_\psi} \left\| \hat{\vartheta}_t^{(s,a)} - \vartheta_P^{(s,a)} \right\|_{\mathcal{H}_\psi} \\ &= 1 + \|1\|_{\mathcal{H}_\psi} \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) . \end{aligned} \quad (21)$$

Combining Equation 20 and Equation 21 with Equation 19 yields:

$$\begin{aligned} & \left| \alpha_t(s, a)^\top v_{h+1}^t - \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)] \right| \\ &\leq B_V \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) + \zeta \left( 2 + \|1\|_{\mathcal{H}_\psi} \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) \right) \\ &\leq \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) + 2\zeta . \end{aligned} \quad (22)$$

Finally, the result follows by noting that

$$\begin{aligned} & Q_h^t(s, a) - (R(s, a) + \mathbb{E}_{X \sim P(\cdot|s,a)} [V_{h+1}^t(X)]) \\ &\leq \alpha_t(s, a)^\top v_{h+1}^t - \mathbb{E}_{P(\cdot|s,a)}[V_{h+1}^t] + \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) \\ &\leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi,t}(s, a) + 2\zeta , \end{aligned}$$

which concludes the proof.  $\blacksquare$

With the result above, we can derive an upper bound on the optimal value  $Q_h^*$  as follows.

**Lemma 7** *Let  $\Theta_P \in \mathcal{C}_t$ . Then  $\forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A}, Q_h^*(s, a) \leq Q_h^t(s, a) + 2(H-h)\zeta$ .*

**Proof** We prove the lemma by induction on  $h$ . When  $h = H$ , the result holds by definition. Now assume that it holds for some  $h' = h + 1$ , where  $1 \leq h < H$ . This implies that

$$\begin{aligned} \forall s \in \mathcal{S}, \quad V_{h+1}^t(s) &= \min \left\{ H, \max_{a \in \mathcal{A}} Q_{h+1}^t(s, a) \right\} \\ &\geq \min \left\{ H, \max_{a \in \mathcal{A}} Q_{h+1}^*(s, a) \right\} - 2(H-h-1)\zeta = V_{h+1}^*(s) - 2(H-h-1)\zeta. \end{aligned}$$

We then have, for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , that

$$\begin{aligned} Q_h^*(s, a) &= R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^*(X)] \\ &\leq R(s, a) + \mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^t(X)] + 2(H-h-1)\zeta. \end{aligned} \quad (23)$$

Using Equation 22 in the proof of Lemma 6, we now see that

$$\mathbb{E}_{X \sim P(\cdot | s, a)} [V_{h+1}^t(X)] \leq \alpha_t(s, a)^\top v_{h+1}^t + \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi, t}(s, a) + 2\zeta$$

which holds as  $\Theta_P \in \mathcal{C}_t$ . We then have from Equation 23 that

$$\begin{aligned} Q_h^*(s, a) &\leq R(s, a) + \alpha_t(s, a)^\top v_{h+1}^t + \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi, t}(s, a) + 2(H-h)\zeta \\ &= Q_h^t(s, a) + 2(H-h)\zeta, \end{aligned}$$

which follows from the definition of  $Q$ -estimates. ■

Given Lemma 6 and Lemma 7, we can finally prove Theorem 3.

### B.3.1. PROOF OF THEOREM 3

If  $\Theta_P \in \mathcal{C}_t$  for all  $t \geq 1$ , we have from Lemma 7 that  $V_1^*(s_1^t) \leq V_1^t(s_1^t) + 2(H-1)\zeta$ . Then following similar steps as in the proof of Theorem 2 and Lemma 3, with probability at least  $1 - \delta/2$ , we have the following:

$$\begin{aligned} \mathcal{R}(N) &\leq \sum_{t \leq T, h \leq H} Q_h^t(s_h^t, a_h^t) - \left( R(s_h^t, a_h^t) + \mathbb{E}_{X \sim P(\cdot | s_h^t, a_h^t)} [V_{h+1}^t(X)] \right) \\ &\quad + 2H \sqrt{2N \log(2/\delta)} + 2\zeta T(H-1) \\ &\leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \sum_{t \leq T, h \leq H} \lambda^{-1/2} \beta_t(\delta/2) \sigma_{\varphi, t}(s_h^t, a_h^t) + 2\zeta TH + 2H \sqrt{2N \log(2/\delta)} \\ &\quad + 2\zeta T(H-1) \\ &\leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \beta_T(\delta/2) \sum_{t \leq T, h \leq H} \lambda^{-1/2} \sigma_{\varphi, t}(s_h^t, a_h^t) + 4\zeta TH + 2H \sqrt{2N \log(2/\delta)} \\ &\leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \alpha_{N, \delta} \sqrt{TH \sum_{t \leq T, h \leq H} \lambda^{-1} \sigma_{\varphi, t}^2(s_h^t, a_h^t) + 4\zeta TH + 2H \sqrt{2N \log(2/\delta)}} \\ &\leq 2 \left( B_V + \zeta \|1\|_{\mathcal{H}_\psi} \right) \alpha_{N, \delta} \sqrt{2(1 + \lambda^{-1} B_\varphi^2 H) N \gamma_N + 4\zeta N + 2H \sqrt{2N \log(2/\delta)}}. \end{aligned}$$

The second step follows from Lemma 6, the third step from monotonicity of  $\beta_t(\delta)$  with  $t$ , the fourth step is due to Cauchy-Schwartz's inequality and the fact that  $\beta_T(\delta/2) \leq \alpha_{N, \delta}$ ,

and the final step follows from [Lemma 5](#). The proof now can be completed using [Theorem 1](#) and taking a union bound.